

# Spiking Neural Networks for Temporal Processing: Status Quo and Future Prospects

Chenxiang Ma\*, Xinyi Chen\*, Yanchen Li\*, Qu Yang, Yujie Wu, Guoqi Li, *Member, IEEE*, Gang Pan, *Senior Member, IEEE*, Huajin Tang, *Senior Member, IEEE*, Kay Chen Tan, *Fellow, IEEE*, Jibin Wu, *Member, IEEE*

**Abstract**—Temporal processing is fundamental for both biological and artificial intelligence systems, as it enables the comprehension of dynamic environments and facilitates timely responses. Spiking Neural Networks (SNNs) excel in handling such data with high efficiency, owing to their rich neuronal dynamics and sparse activity patterns. Given the recent surge in the development of SNNs, there is an urgent need for a comprehensive evaluation of their temporal processing capabilities. In this paper, we first conduct an in-depth assessment of commonly used neuromorphic benchmarks, revealing critical limitations in their ability to evaluate the temporal processing capabilities of SNNs. To bridge this gap, we further introduce a benchmark suite consisting of three temporal processing tasks characterized by rich temporal dynamics across multiple timescales. Utilizing this benchmark suite, we perform a thorough evaluation of recently introduced SNN approaches to elucidate the current status of SNNs in temporal processing. Our findings indicate significant advancements in recently developed spiking neuron models and neural architectures regarding their temporal processing capabilities, while also highlighting a performance gap in handling long-range dependencies when compared to state-of-the-art non-spiking models. Finally, we discuss the key challenges and outline potential avenues for future research.

**Index Terms**—Spiking Neural Networks, Temporal Processing, Neuromorphic Benchmarks, Neuromorphic Computing

## I. INTRODUCTION

TEMPORAL processing is a fundamental capability that enables animals to perceive, plan, and act within dynamic environments. To effectively analyze and understand temporal data, a variety of models have been developed based on artificial neural networks (ANNs), such as Recurrent Neural Networks (RNNs) [1], [2], Temporal Convolutional Networks (TCNs) [3], Transformers [4], and State Space Models (SSMs) [5]. Despite their impressive performance in processing temporal data, these models often demand significant

computational resources, which can limit their deployment on resource-constrained platforms [6].

In contrast, Spiking Neural Networks (SNNs) [7], inspired by the computational principles of biological neural networks, offer an energy-efficient computational framework for temporal processing [8]. By employing a sparse spike-based representation, SNNs enable efficient event-driven computation, wherein spiking neurons activate solely in response to incoming spikes [9], [10]. This characteristic is particularly advantageous when implemented on neuromorphic chips [11]–[15], where SNNs have demonstrated significantly enhanced energy efficiency compared to traditional ANNs [16]. Beyond their energy efficiency, spiking neurons inherently function as stateful models, characterized by rich neuronal dynamics that arise from complex morphology, variations in ionic conductance, and the distribution of synaptic inputs [17]. These features endow SNNs with significant potential for representing and processing temporal data [18]–[20].

Recently, numerous approaches have been developed to enhance training efficiency [21]–[30] and representation power of SNNs [31]–[37]. Despite these advancements, a clear consensus on how to evaluate these approaches in the context of temporal processing remains elusive. Additionally, the lack of standardized training and evaluation configurations across these studies complicates fair comparisons. These challenges collectively impede the progress of the field and limit the real-world applicability of SNNs. In this paper, we seek to address these issues by establishing a comprehensive evaluation benchmark specifically focused on temporal processing capabilities. Furthermore, we will assess existing approaches to elucidate the current state of the field and identify potential future research directions. An overview of the paper's structure is presented in Fig. 1.

Benchmarks are essential in Artificial Intelligence research, as they provide standardized datasets and evaluation metrics that facilitate consistent performance comparison, track progress, and promote reproducibility. Currently, the benchmarks commonly employed for SNN evaluation can be categorized into three groups. The first category comprises static image recognition datasets [38], [39], which require the conversion of static images into sequences, typically by replicating images along the time axis [40]. The second category encompasses event-based vision datasets generated using neuromorphic sensors, such as Dynamic Vision Sensor (DVS) cameras. Some of these datasets are created by imposing artificial saccadic motion on static images [41], [42], while others directly capture real-world moving objects

\*Chenxiang Ma, Xinyi Chen, and Yanchen Li contributed equally to this article. Corresponding Author: Jibin Wu (jibin.wu@polyu.edu.hk)

Chenxiang Ma, Xinyi Chen, Yanchen Li, Kay Chen Tan, and Jibin Wu are with the Department of Data Science and Artificial Intelligence, The Hong Kong Polytechnic University, Hong Kong SAR. Jibin Wu is also with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong SAR

Qu Yang is with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 119077

Yujie Wu is with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong SAR

Guoqi Li is with the Institute of Automation, Chinese Academy of Sciences, Beijing 100045, China

Gang Pan and Huajin Tang are with the State Key Laboratory of Brain-Machine Intelligence, College of Computer Science and Technology, MOE Frontier Science Center for Brain Science and Brain-Machine Integration, Zhejiang University, Hangzhou 310027, China

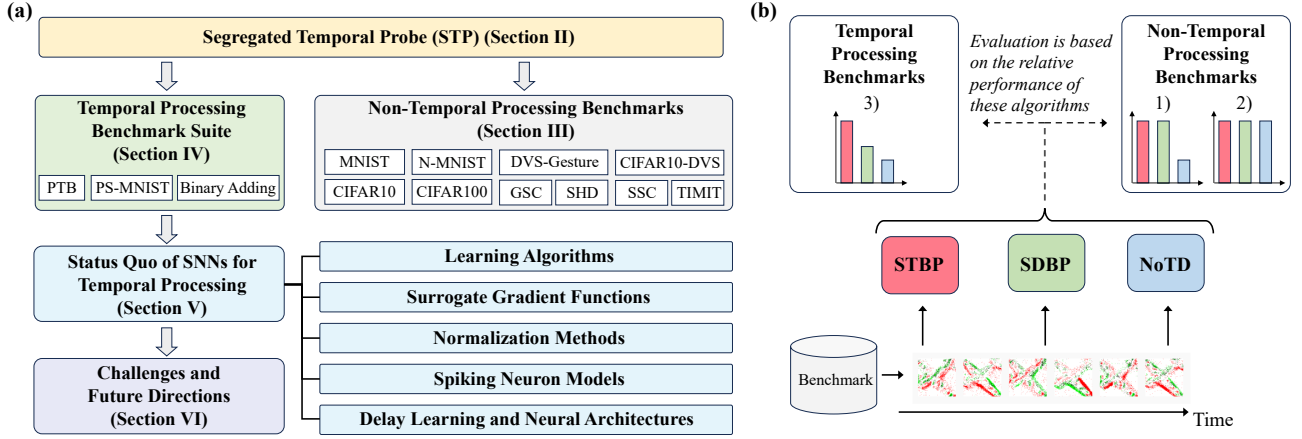


Fig. 1. (a) Overview of the paper organization. In Section II, we propose the Segregated Temporal Probe (STP) analytical tool for assessing the effectiveness of neuromorphic benchmarks in evaluating the temporal processing capabilities of SNN. In Section III, using the STP, we discover that commonly used neuromorphic benchmarks are ineffective for assessing temporal processing performance. Then, we introduce a suite of temporal processing benchmarks in Section IV. Based on this suite, we conduct a comprehensive benchmarking study to reveal the current status of SNNs for temporal processing in Section V. Finally, we discuss key challenges and outline future directions in Section VI. (b) Illustration of the proposed STP. STP incorporates three algorithms (STBP, SDBP, and NoTD) that systematically disrupts the temporal processing pathways within an SNN to elucidate their significance.

using DVS cameras [43], [44]. The third category consists of audio classification datasets, which are either created using the original audio signals [45], [46] or converted into spike-based representations using biologically inspired encoding methods [47], [48]. While these benchmarks have significantly advanced neuromorphic computing research over the past decade [8], [40], their effectiveness in evaluating the temporal processing capability remains unclear, potentially leading to inaccurate assessments and misleading conclusions.

To address this issue, we propose an analytical tool called the Segregated Temporal Probe (STP), designed to evaluate the effectiveness of a benchmark dataset in assessing the temporal processing capabilities of SNNs. Specifically, STP incorporates three learning algorithms: Spatio-Temporal Backpropagation (STBP) [49], Spatial Domain Backpropagation (SDBP), and No Temporal Domain (NoTD). These algorithms systematically disrupt the temporal processing pathways within an SNN to elucidate their significance. In particular, STBP preserves the full temporal processing pathways of the SNN during both forward and backward passes, whereas SDBP stops gradient propagation along the temporal domain during the backward pass. In contrast, NoTD stops the propagation of information along the temporal domain during both passes, effectively treating each time step independently.

By applying STP to widely adopted neuromorphic benchmark datasets, we find that NoTD achieves performance comparable to STBP on static image recognition [38], [39] and event-based vision datasets [41]–[43]. This finding suggests that these datasets can be effectively processed without relying on temporal processing capabilities of SNN models. For audio classification datasets [45], [46], [48], both SDBP and STBP obtain similar performance, suggesting that temporal credit assignment during the backward pass is not necessary for these datasets. Consequently, these existing benchmark datasets do not adequately evaluate the temporal processing capability of SNNs.

To bridge this gap and elucidate the status quo of the temporal processing capabilities of existing SNN approaches, we introduce a benchmark suite encompassing three temporal processing tasks characterized by rich temporal dynamics. Subsequently, we conduct a comprehensive evaluation of over thirty SNN approaches using this benchmark suite. Our benchmarking study reveals three significant findings that have not been previously reported: (1) Online learning algorithms [50], [51], which claim performance comparable to or even superior to STBP, often yield less competitive results on temporal processing tasks. This indicates that the omitted temporal gradients are crucial for learning temporal dependencies within such data. (2) Surrogate gradient functions [52] with smoother curves and reduced gradient mismatching, such as Triangle [25] and Sigmoid [53], prove to be more effective for temporal processing tasks. (3) Recent advancements in spiking neuron models that incorporate strategies such as additional memory states [54], [55], heterogeneous neuron parameters [56], [57], and enriched recurrent neuron dynamics [58]–[60] demonstrate significant improvements over the simplified Leaky Integrated-and-Fire (LIF) model [61] in temporal processing. However, despite their considerable energy efficiency, these state-of-the-art (SOTA) SNN models still lag behind ANN models [2], [5] in their ability to model long-range temporal dependencies.

Our major contributions in this work are summarized as follows:

- We propose an analytical tool, STP, for evaluating the effectiveness of neuromorphic benchmarks in assessing temporal processing capabilities. This tool facilitates the development of neuromorphic benchmarks specifically tailored for temporal processing.
- We identify critical limitations in existing neuromorphic benchmarks regarding their evaluation of temporal processing capabilities and propose a new benchmark suite explicitly designed for this purpose.

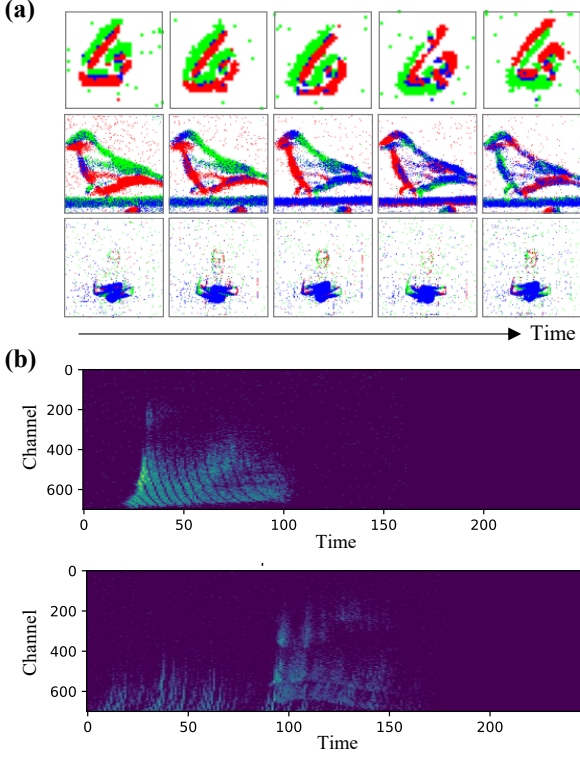


Fig. 2. Visualization of samples in neuromorphic benchmarks. (a) Samples from event-based vision datasets: N-MNIST, CIFAR10-DVS, and DvsGesture (from top to bottom). (b) Samples from neuromorphic audio datasets: SHD (top) and SSC (bottom).

- We conduct a comprehensive benchmarking study of over thirty SNN approaches to elucidate the current state of the field.
- We develop an open-sourced library<sup>1</sup> for neuromorphic temporal processing, which enables consistent performance comparisons across different approaches and facilitates the tracking of advancements in the field.

The remainder of this paper is organized as follows. Section II introduces the proposed STP tool, which is utilized to evaluate existing neuromorphic benchmarks for temporal processing in Section III. Section IV presents a novel neuromorphic benchmark suite specifically designed for temporal processing. Subsequently, we conduct a comprehensive evaluation of over thirty SNN approaches in Section V to elucidate the current state of the field. Section VI discusses key challenges and outlines potential future research directions. Finally, we conclude the paper in Section VII.

## II. SEGREGATED TEMPORAL PROBE (STP)

While existing neuromorphic benchmarks have significantly advanced the field of neuromorphic computing [8], [40], it remains unclear whether these benchmarks effectively evaluate critical temporal information. This uncertainty is particularly evident in benchmarks adapted from static datasets [41], [42], where objects in individual frames often provide sufficient

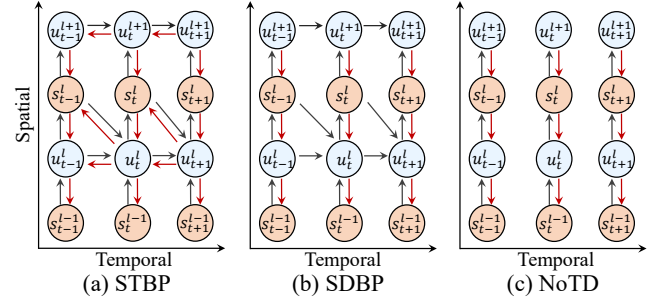


Fig. 3. Comparison of the computational graphs for three algorithms utilized in the STP. Forward and backward passes are denoted by black and red arrows, respectively.

information for classification, as illustrated in Fig. 2(a). Consequently, we are motivated to explicitly analyze whether temporal processing capability is genuinely critical for achieving high performance on these benchmarks. To this end, we introduce the STP, which disrupts the temporal processing pathways within an SNN to elucidate their significance.

As illustrated in Fig. 1(b), the STP tool comprises three learning algorithms: STBP [49], SDBP, and NoTD. STBP serves as the baseline, retaining the entire temporal processing pathways, including both the forward pass of activation values and the backward pass of error gradients across time. In contrast, SDBP is designed to disrupt temporal processing during the backward pass while preserving it in the forward pass. NoTD, on the other hand, eliminates temporal processing entirely by processing each frame independently at each time step. In the following, we will use the LIF neuron model as an example to illustrate these algorithms. Fig. 3 provides a visualization that highlights the differences in the forward and backward passes of each algorithm.

*a) LIF Neuron Model:* The LIF neuron model [61] is one of the most widely used spiking neuron models due to its simplicity and analytical tractability [8]. As described in Eqs. (1)–(3), LIF neurons capture the dynamics of membrane potential, which continuously integrates input spikes from preceding neurons. When the membrane potential exceeds a specified firing threshold, a spike is generated, followed by a reset of the membrane potential to its resting state.

$$\mathbf{u}^l[t] = \underbrace{\lambda \cdot \mathbf{u}^l[t-1] \cdot (1 - s^l[t-1])}_{\text{temporal processing}} + \mathbf{W}^l \cdot \mathbf{s}^{l-1}[t], \quad (1)$$

$$s^l[t] = \Theta(\mathbf{u}^l[t] - V_{\text{th}}), \quad (2)$$

$$\Theta(x) = \begin{cases} 1, & x \geq 0, \\ 0, & \text{otherwise,} \end{cases} \quad (3)$$

where  $\mathbf{u}^l[t]$  represents the membrane potential in layer  $l$  at time step  $t$ ,  $\lambda$  denotes a decay factor that determines the rate at which the membrane potential decays over time, and  $\mathbf{s}^{l-1}[t]$  represents input spikes from the previous layer.  $\mathbf{W}^l$  is the weight matrix,  $\Theta(x)$  is the Heaviside step function, and  $V_{\text{th}}$  denotes the firing threshold. The resting potential is typically set to zero in practice.

*b) STBP:* STBP is a gradient-based learning algorithm specifically designed for SNNs [49]. After input spikes are

<sup>1</sup>Code is publicly available at <https://github.com/liyc5929/neuroseqbench>.

propagated to the output layer  $L$ , a loss  $\mathcal{L}$  is calculated by comparing the predictions with the target values. The backward pass then initiates, calculating the gradient of the loss with respect to each model parameter. By employing the chain rule, the gradient of the loss with respect to the weights of layer  $l$  can be computed as follows:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} = \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t]} \cdot \frac{\partial \mathbf{u}^l[t]}{\partial \mathbf{W}^l} = \sum_{t=1}^T \delta^l[t]^\top \cdot \mathbf{s}^{l-1}[t]^\top, \quad (4)$$

$$\delta^l[t] = \begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{u}^L[T]}, & l = L \text{ and } t = T, \\ \delta^L[t+1] \cdot \frac{\partial \mathbf{u}^L[t+1]}{\partial \mathbf{u}^L[t]} + \frac{\partial \mathcal{L}}{\partial \mathbf{u}^L[t]}, & l = L \text{ and } t < T, \\ \delta^{l+1}[T] \cdot \frac{\partial \mathbf{u}^{l+1}[T]}{\partial \mathbf{s}^l[T]} \cdot \frac{\partial \mathbf{s}^l[T]}{\partial \mathbf{u}^l[t]}, & l < L \text{ and } t = T, \\ \delta^l[t+1] \frac{\partial \mathbf{u}^{l+1}[t+1]}{\partial \mathbf{u}^l[t]} + \delta^{l+1}[t] \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{u}^l[t]}, & \text{otherwise,} \end{cases} \quad (5)$$

where  $\delta^l[t] \triangleq \frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t]}$ , represents the error back-propagated from the last layer and the last time step.  $T$  denotes the total number of time steps. The Heaviside step function is non-differentiable since its gradient, i.e.,  $\frac{\partial s_i^l[t]}{\partial u_i^l[t]}$ , is zero except at the point where  $\mathbf{u}_i^l[t] = V_{th}$ , where it becomes infinite. To address this issue, a continuous surrogate gradient function [52] is often adopted to replace the non-differentiable Heaviside step function during the backward pass, as denoted by  $\frac{\partial s_i^l[t]}{\partial u_i^l[t]} \approx \mathbb{H}(\mathbf{u}_i^l[t])$ .

c) *SDBP*: To illustrate SDBP, we begin by explicitly detailing the gradients associated with the temporal processing function. This is achieved by rewriting the error term  $\delta^l[t]$  under the general condition ( $l < L$  and  $t < T$ ):

$$\delta^l[t] = \delta^{l+1}[t] \cdot \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \cdot \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} + \underbrace{\sum_{t'=t+1}^T \delta^{l+1}[t'] \cdot \frac{\partial \mathbf{u}^{l+1}[t']}{\partial \mathbf{u}^l[t']}}_{\text{temporal processing gradients}} \cdot \frac{\partial \mathbf{u}^l[t']}{\partial \mathbf{u}^l[t]}, \quad (6)$$

where the gradient at the current time step  $t$  is influenced by a sum of the gradients from all subsequent time steps. This stems from the temporal processing function in Eq. (1). More concretely, the membrane potential at the current time step  $t$  implicitly contributes to all membrane potentials at the subsequent time steps due to the decay and resetting processes over time. Accordingly, the errors at all subsequent steps need to be included in the computation of the error in the time  $t$ .

SDBP keeps the same temporal processing function in the forward pass as STBP but omits it in the backward pass. Consequently, error signals generated at a given time step cannot be propagated back to previous time steps, limiting the use of error information to refine earlier predictions. Formally, let  $\epsilon^l[t]$  represent  $\frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t]}$  in SDBP, the weight gradient at the layer  $l$  is given by

$$\nabla_{\mathbf{W}^l} \mathcal{L} = \sum_{t=1}^T \epsilon^l[t]^\top \cdot \frac{\partial \mathbf{u}^l[t]}{\partial \mathbf{W}^l} = \sum_{t=1}^T \epsilon^l[t]^\top \cdot \mathbf{s}^{l-1}[t]^\top, \quad (7)$$

TABLE I  
EVALUATION RESULTS OF WIDELY-USED NEUROMORPHIC BENCHMARKS. “Acc.” STANDS FOR “ACCURACY”.

Dataset	Time Step ( $T$ )	Method	Acc.	$\Delta$ Acc.
MNIST [38]	10	STBP	99.40	-
		SDBP	99.27	-0.13
		NoTD	99.18	-0.22
CIFAR10 [39]	4	STBP	94.86	-
		SDBP	94.74	-0.12
		NoTD	93.46	-1.40
CIFAR100 [39]	4	STBP	74.57	-
		SDBP	74.35	-0.22
		NoTD	73.28	-1.29
N-MNIST [41]	300	STBP	99.49	-
		SDBP	99.48	-0.01
		NoTD	99.09	-0.40
CIFAR10-DVS [42]	10	STBP	78.50	-
		SDBP	79.00	+0.50
		NoTD	80.00	+1.50
DvsGesture [43]	20	STBP	95.14	-
		SDBP	95.83	+0.69
		NoTD	94.44	-0.70
GSC [45]	101	STBP	92.91	-
		SDBP	89.00	-3.91
		NoTD	77.53	-15.38
SHD [48]	250	STBP	86.48	-
		SDBP	85.07	-1.41
		NoTD	68.51	-17.97
SSC [48]	250	STBP	67.13	-
		SDBP	66.03	-1.10
		NoTD	44.97	-22.16
TIMIT [46]	100	STBP	57.07	-
		SDBP	53.24	-3.83
		NoTD	49.01	-8.06

$$\epsilon^l[t] = \begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{u}^L[t]}, & l = L, \\ \epsilon^{l+1}[t] \cdot \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{u}^l[t]}, & l < L. \end{cases} \quad (8)$$

d) *NoTD*: NoTD is designed to eliminate temporal processing functions in both forward and backward passes, allowing each time step in the input sequence to be processed independently without interaction with other time steps. Consequently, NoTD cannot learn temporal relationships within sequences or integrate information over time for decision-making. Formally, the forward pass of layer  $l$  in NoTD is given by

$$\mathbf{u}^l[t] = \mathbf{W}^l \cdot \mathbf{s}^{l-1}[t], \quad (9)$$

$$\mathbf{s}^l[t] = \Theta(\mathbf{u}^l[t] - V_{th}). \quad (10)$$

Compared to Eq. (1), Eq. (9) removes the temporal processing function, i.e., the membrane potential update process. In the backward pass, the weight gradient in NoTD is the same as Eqs. (7) and (8) in SDBP.

e) *Evaluation Criteria*: To determine whether a dataset can effectively evaluate the temporal processing capabilities of SNNs, we train SNNs using the three algorithms on the dataset. The effectiveness of the dataset will be evaluated based on the relative performance of these algorithms:

- 1) If the performance of SDBP is comparable to or exceeds that of STBP, this suggests that temporal credit assignment during the backward pass may be unnecessary, indicating that the dataset is not suitable for evaluating temporal processing capabilities.

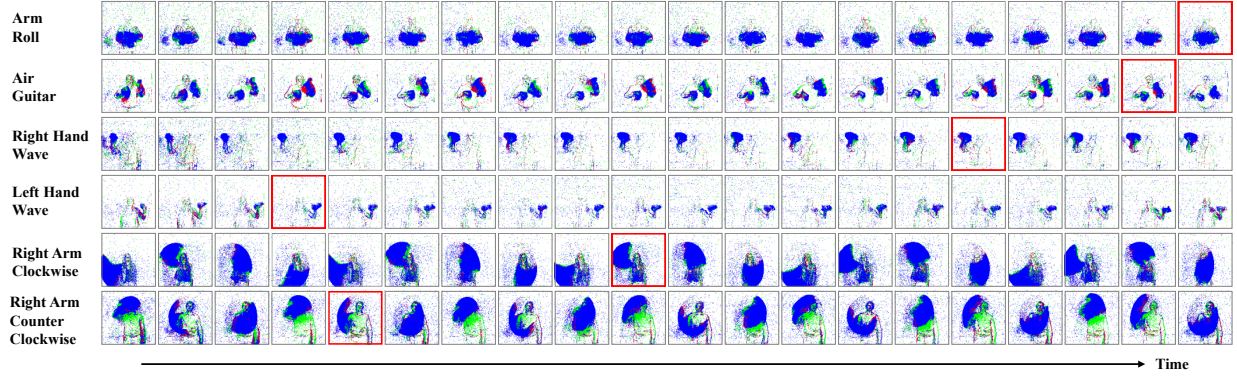


Fig. 4. Qualitative results of samples from the DvsGesture dataset, along with the confident frame (highlighted with red boxes) selected by all the algorithms in STP. Labels are provided in the leftmost column.

- 2) If the performance of NoTD is comparable to or better than that of SDBP, it implies that the dataset can be effectively addressed without incorporating temporal interactions within the model, thus rendering it an inadequate benchmark.
- 3) Conversely, if SDBP surpasses NoTD and STBP outperforms SDBP, this indicates that information integration across the time is essential, and that temporal credit assignment effectively facilitate this integration. Therefore, the benchmark is suitable for assessing the temporal processing capabilities of SNNs.

### III. EVALUATION OF NEUROMORPHIC BENCHMARKS USING STP

We conduct a comprehensive analysis of the existing neuromorphic benchmarks using the STP to reveal their critical limitations. For each benchmark, we adhere to standard protocols from prior studies [49]–[51], [54], [62], [63], which include dataset processing, data augmentation, network design, and training configurations. Detailed descriptions of these setups are provided in Appendix A.

Our analysis begins with static image recognition benchmarks, including MNIST [38], CIFAR10 [39], and CIFAR100 [39]. These datasets do not contain temporal information, as each input sequence is generated by repeating a static image along the time dimension. This is reaffirmed by our experimental results. As shown in Table I, both NoTD and SDBP achieve accuracy comparable to STBP. This suggests that the capability to model temporal relationships is not necessary for high performance on these datasets.

Next, we examine event-based vision datasets. Results are presented in Table I. For synthetic event-based datasets like N-MNIST [41] and CIFAR10-DVS [42], NoTD performs similarly to or even outperforms both SDBP and STBP. This suggests that these datasets can be effectively addressed at the frame level without the need for temporal modeling. Interestingly, the DvsGesture dataset [43], which captures human gestures in real time, is also effectively addressed by NoTD. This implies that the DvsGesture dataset does not require the model to have temporal processing capabilities for accurate recognition.

To understand these results, we conduct a qualitative analysis by visualizing samples from the DvsGesture dataset in Fig. 4, along with the confident frame selected by each algorithm in

STP. The confident frame is defined as the frame where the output layer has the strongest response. We observe that most samples have minimal changes over time, which allows accurate classification based on a single frame. The three algorithms often select the same confident frame, reinforcing that a single frame is sufficient for correct recognition without the need for temporal modeling.

In cases where the algorithms choose different frames, as shown in Appendix Fig. 9, the selected frames often have similar spatial features despite at different time steps. This demonstrates that these samples contain several informative frames adequate for pattern recognition. Additionally, certain classes, such as right arm clockwise and counterclockwise movements, seemingly necessitate temporal modeling due to their directional differences. However, both NoTD and SDBP successfully recognize these classes and select confident frames that align with STBP. This effectiveness can be attributed to the presence of distinguishable spatial features within these classes, which reduces the necessity for temporal modeling. Furthermore, we visualize samples correctly classified by STBP but not by NoTD in Appendix Fig. 10. We observe that the errors made by NoTD are primarily associated with spatial features rather than temporal cues. Similarly, as illustrated in Appendix Fig. 11, the misclassifications by STBP are also linked to spatial features present within individual frames, rather than to temporal relationships across multiple frames. This analysis supports the conclusion that the primary challenge for achieving high-accuracy classification on the DvsGesture dataset resides in spatial pattern recognition rather than in temporal modeling.

Finally, we evaluate audio classification benchmarks, including GSC [45], SHD [48], SSC [48], and TIMIT [46]. Fig. 2(b) presents a visualization of randomly selected audio samples. By applying STP to these benchmarks, we find that a significant subset of samples can be accurately classified using only frame-level processing capabilities. This is evidenced by the moderate accuracy achieved by NoTD on these benchmarks. Furthermore, SDBP demonstrates substantially higher accuracy than NoTD, approaching the performance of STBP. This indicates that temporal interactions during the forward pass are critical for these tasks, while temporal credit assignment during the backward pass contributes only marginally to performance gains. Consequently, these benchmarks are not effective for

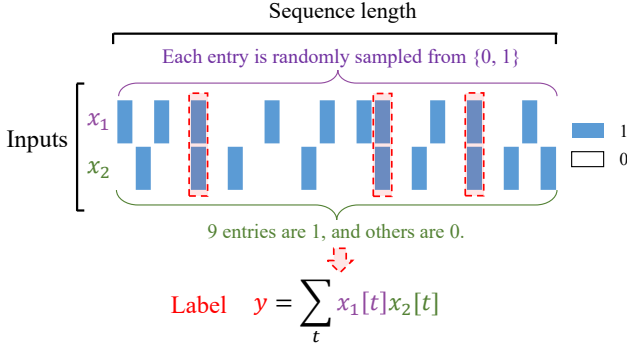


Fig. 5. Illustration of the binary adding task, designed to test the ability of SNN models to capture long-range dependencies. The sequence length in this task is adjustable, allowing flexibility in controlling the task’s difficulty.

evaluating temporal processing capabilities.

In summary, our analysis utilizing the STP tool reveals that the current neuromorphic benchmarks are inadequate for evaluating the temporal processing capabilities of SNNs.

#### IV. TEMPORAL PROCESSING BENCHMARK SUITE

To bridge this gap, we present a new benchmark suite comprising three tasks in this section, along with a validation study to demonstrate their effectiveness in assessing temporal processing capabilities.

##### A. Benchmark Suite

To reveal the status quo of SNNs for temporal processing, we present a benchmark suite designed to comprehensively evaluate the temporal processing capabilities of existing SNN approaches. This suite incorporates three tasks with distinct modalities: language generation, pixel-level image classification, and mathematics.

- **Penn Treebank (PTB).** PTB is a widely used language modeling dataset [64], derived from Wall Street Journal articles, including various text types such as news reports, editorials, and financial analyses. In this task, the text is tokenized into individual words with a vocabulary size of 10,000. Each sample is then truncated into sequences of 70 tokens. The model is required to predict the next token at each time step based on the preceding input context.
- **Permuted-Sequential MNIST (PS-MNIST).** PS-MNIST is a sequence classification dataset derived from the MNIST dataset [38]. Each gray-scaled image from the MNIST dataset is first flattened in row-major order into a sequence of length 784. The pixel order is then shuffled using a fixed random permutation matrix, disrupting locality to challenge the model’s ability to model long-range dependencies. Finally, these pixels are fed to an SNN one pixel at a time, with the prediction made at the final time step.
- **Binary Adding.** To challenge models in performing temporal processing over long distances, we propose a novel binary adding task. As illustrated in Fig. 5, we generate a synthetic 10-class dataset with two input channels, denoted as  $\{x_1, x_2\} \in \{0, 1\}^T$ , where  $x_1$  represents a binary value

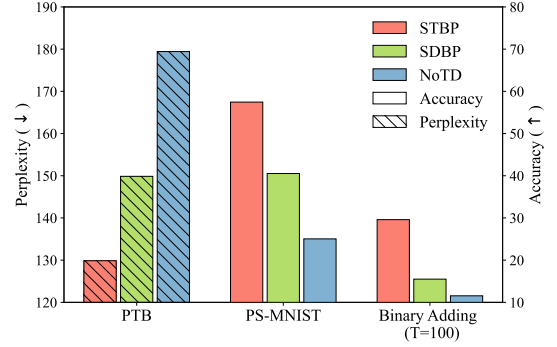


Fig. 6. Validation of the temporal processing benchmark suite through the STP

sequence and  $x_2$  serves as an index indicator. The input sequence length  $T$  can be flexibly adjusted to change the task difficulty. The binary value sequences  $x_1$  consist of  $T$  entries randomly sampled from  $\{0, 1\}$ , while the indicator sequence  $x_2$  assigns a value of 1 to nine randomly selected indices and 0 otherwise. The indicator sequence  $x_2$  acts as a pointer to indicate which entries in  $x_1$  should be added. The label is defined as  $y = \sum_{t=1}^T x_1[t] \cdot x_2[t]$ , which ranges from 0 to 9. The model must process the entire sequence before making a prediction, necessitating the ability to integrate information over a long time span. For a consistent comparison across different SNN approaches, we construct a fixed dataset containing 50,000 training samples and 2,000 testing samples.

##### B. Validation of Benchmarks using the STP

We further apply the STP tool to validate the effectiveness of these three benchmarks in assessing temporal processing capabilities. Details of the training setups are provided in Appendix B. As shown in Fig. 6, STBP significantly outperforms SDBP, which in turn substantially surpasses NoTD. This observation is consistent across all three benchmarks, demonstrating that these benchmarks contain critical temporal information that is necessary to be captured for high performance. Therefore, they can serve as effective benchmarks for evaluating the temporal processing capabilities of various SNN approaches.

#### V. STATUS QUO OF SNNs FOR TEMPORAL PROCESSING

To assess the current state of SNNs in temporal processing, we further conduct a comprehensive study of over thirty recently developed SNN methods using our proposed benchmark suite. The evaluated methods encompass a wide range of aspects, including learning algorithms in Section V-A, surrogate gradients in Section V-B, normalization techniques in Section V-C, neuron models in Section V-D, and neural architectures in Section V-F. For completeness, each method is evaluated using both feedforward and recurrent architectures [65], [66]. The Recurrent SNNs (RSNNs) incorporate recurrent weights to retain historical information, thereby offering improved memory capacity compared to feedforward architectures.

To ensure a fair comparison across these methods, we re-implement each method to the best of our ability, utilizing their publicly available source codes and descriptions provided in

TABLE II

COMPARISON OF LEARNING ALGORITHMS ON TEMPORAL PROCESSING BENCHMARKS. “FF” AND “REC.” REFER TO “FEEDFORWARD” AND “RECURRENT” ARCHITECTURES, RESPECTIVELY. “PPL” STANDS FOR “PERPLEXITY.”

Dataset	PTB ( $T = 70$ )		PS-MNIST ( $T = 784$ )		Binary Adding ( $T = 100$ )	
Metric	PPL ↓		Acc. ↑		Acc. ↑	
Method	FF	Rec.	FF	Rec.	FF	Rec.
STBP [49]	129.96	111.96	57.45	72.97	29.60	53.35
T-STBP	137.8	120.58	53.00	71.03	23.00	51.50
E-prop [66]	-	125.54	-	52.88	-	50.85
OTTT [50]	141.77	-	44.61	-	17.20	-
SLTT [51]	149.86	-	40.53	-	15.50	-

their respective papers. We maintain consistent training setups across all methods. To optimize the hyperparameters for each method, we conduct a grid search; detailed information on the hyperparameters and training setups is available in our open-sourced library. We have open-sourced our code and benchmarks to facilitate future advancements in the field, and we welcome contributions from the neuromorphic computing community to expand our benchmark by including more tasks and recent SNN methods.

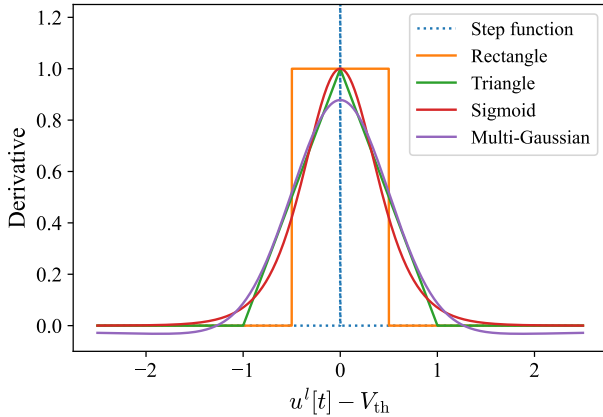


Fig. 7. Comparison of different surrogate gradient functions.

#### A. Learning Algorithms

We begin by evaluating the performance of five SNN learning algorithms, including two offline learning algorithms (i.e., STBP [49], Truncated-STBP (T-STBP)) as well as three recently-proposed online learning algorithms (i.e., E-prop [66], Online Training Through Time (OTTT) [50], and Spatial Learning Through Time (SLTT) [51]). STBP is a standard gradient-based learning algorithm, which unfolds SNNs over both spatial and temporal domains and applies gradient descent to the entire computational graph. T-STBP processes input sequences in smaller segments, preventing error gradients from being backpropagated between these segments while preserving the complete temporal gradients within each individual segment. In contrast, the three online learning algorithms update SNNs at every time step, and they either partially or fully disregard temporal gradients to enable online learning. Specifically, E-prop achieves online learning in RSNNs by using eligibility traces, which accumulate presynaptic activities over time.

TABLE III

COMPARISON OF SURROGATE GRADIENT FUNCTIONS ON TEMPORAL PROCESSING BENCHMARKS. THE BEST MODEL IS HIGHLIGHTED IN **BOLD**, AND THE SECOND BEST IS UNDERLINED.

Dataset	PTB ( $T = 70$ )		PS-MNIST ( $T = 784$ )		Binary Adding ( $T = 100$ )	
Metric	PPL ↓		Acc. ↑		Acc. ↑	
Method	FF	Rec.	FF	Rec.	FF	Rec.
Rectangle [49]	129.96	111.96	57.45	72.97	<u>29.60</u>	53.35
Triangle [25]	127.78	108.98	62.96	<b>77.15</b>	28.20	62.90
Multi-Gaussian [67]	130.01	112.74	57.01	71.96	<b>32.20</b>	<u>65.60</u>
Sigmoid [53]	<b>127.33</b>	<b>107.92</b>	<b>63.18</b>	76.02	29.35	<b>66.45</b>
ASGL [63]	128.14	113.18	58.32	66.00	26.15	46.70

Similarly, OTTT extends the mechanism of eligibility traces to feedforward SNNs, enabling effective online training for large-scale datasets. SLTT, on the other hand, ignores all temporal gradients, removing the need for additional traces while slightly compromising gradient precision.

In our evaluation, we adhere to the original setups for these algorithms, applying OTTT and SLTT exclusively to feedforward SNNs, while E-prop is used solely for RSNNs. Results presented in Table II indicate that STBP consistently outperforms the other four learning algorithms across all three benchmarks, regardless of the architecture. This performance gap is attributed to the omission of temporal gradients in these algorithms, which leads to biased gradient estimations and impedes the accurate propagation of temporal errors.

Notably, our findings contrast with previous studies [50], [51], [66] that reported nearly lossless performance for these online learning algorithms compared to STBP. This discrepancy stems from the fact that earlier evaluations were primarily on datasets with no or limited temporal information, such as CIFAR10, CIFAR10-DVS, and DvsGesture, where the absence of temporal gradients had a minimal effect. In contrast, our benchmarks necessitate effective learning over time, where the accurate calculation of temporal gradients is crucial. This experiment underscores the limitations of commonly used neuromorphic benchmarks and highlights a significant accuracy gap between existing online learning methods and STBP in temporal processing tasks.

#### B. Surrogate Gradient Functions

To address the nondifferentiable activation functions used in SNNs, surrogate gradient methods [68] are commonly employed. These methods retain the original step function during the forward pass while replacing it with a smooth and continuous function during the backward pass. In this section, we evaluate the performance of five most frequently used surrogate gradient functions: Rectangle [49], [69], Triangle [25], [68], Sigmoid [53], Multi-Gaussian [67], and Adaptive Smoothing Gradient Learning (ASGL) [63]. Notably, ASGL [63] not only smooths the step function during backpropagation but also partially replaces the step function in the forward pass with the integral of surrogate functions during training. This approach mitigates the gradient approximation error caused by the mismatch between actual and surrogate gradients, resulting in smoother training.

Results in Table III indicate that while the performance of different surrogate functions shows only minor differences in

TABLE IV  
COMPARISON OF NORMALIZATION METHODS ON TEMPORAL PROCESSING BENCHMARKS

Dataset	PTB ( $T = 70$ )		PS-MNIST ( $T = 784$ )		Binary Adding ( $T = 100$ )	
Metric	PPL ↓		Acc. ↑		Acc. ↑	
Method	FF	Rec.	FF	Rec.	FF	Rec.
w/o BN	129.96	111.96	57.45	72.97	29.60	53.35
TEBN [25]	126.88	102.24	<b>94.94</b>	<b>95.02</b>	<b>53.03</b>	<b>65.45</b>
TDBN [70]	127.04	101.54	72.74	88.60	35.55	64.10
LayerNorm [71]	<b>123.54</b>	<b>101.53</b>	62.28	69.96	40.25	49.70

feedforward architectures, their impact is significantly more pronounced in recurrent architectures. This is due to the recurrent connection from the output spike at time  $t - 1$  to the membrane potential at time  $t$ , which introduces an additional pathway for gradient backpropagation. This iterative computation of surrogate gradients over time introduces additional challenges for SNN training. Fortunately, a well-shaped surrogate function can more effectively propagate gradients back to earlier time steps through this recurrent pathway, resulting in improved network convergence for recurrent architectures.

To identify the most suitable surrogate gradient functions for temporal processing tasks, we further rank their performance across the three tasks and two architectures. The average ranking is used as a score to quantify their effectiveness. Our results indicate that the Sigmoid [53] and Triangle [25] functions achieve the top two average rankings, with scores of 0.67 and 1.17, respectively. This finding suggests that these two functions are more appropriate for temporal tasks. The superior performance of the Sigmoid and Triangle functions can be attributed to two factors. Firstly, their higher smoothness facilitates error propagation during SNN training [25], [40], [49]. Additionally, the values of these surrogate functions provide a closer approximation to the ill-defined derivative of the step function, as illustrated by the dotted curves in Fig. 7. Consequently, the gradient mismatch problem [63] can be more effectively alleviated, contributing to better training convergence. Finally, we examine the performance of the ASGL strategy [63]. Despite the performance improvements observed in previous static tasks, our results indicate that this strategy is less effective for temporal processing tasks. This inefficacy primarily arises from the discrepancy between the spike generation function employed during training and that used during inference, which accumulates over time, particularly when handling long sequences.

### C. Normalization Methods

Here, we evaluate the influence of normalization methods on the temporal processing performance of SNNs. Building upon the Batch Normalization (BN) [72] method in ANNs, normalization methods have been developed specifically for SNNs to improve training stability and accelerate convergence. For example, Threshold-Dependent BN (TDBN) [70] normalizes features across both batch and time dimensions and adjusts the normalized variance based on the threshold, which effectively mitigates the problems of gradient vanishing or explosion during training. Furthermore, Temporal Effective BN (TEBN) [62]

TABLE V  
COMPARISON OF SPIKING NEURON MODELS ON TEMPORAL PROCESSING BENCHMARKS. THE **BEST** MODEL IS HIGHLIGHTED IN **BOLD**, THE SECOND BEST IS UNDERLINED, AND THE *THIRD* BEST IS *ITALICIZED*.

Dataset	PTB ( $T = 70$ )		PS-MNIST ( $T = 784$ )		Binary Adding ( $T = 100$ )	
Metric	PPL ↓		Acc. ↑		Acc. ↑	
Network	FF	Rec.	FF	Rec.	FF	Rec.
#Params.	~5M	~6M	~90K	~160K	~20K	~40K
LIF	129.96	111.96	57.45	72.97	29.60	53.35
PLIF [56]	123.76	105.64	55.86	77.32	29.40	53.25
ALIF [67]	113.67	102.25	73.90	85.78	40.30	68.00
adLIF [53]	118.52	<b>97.22</b>	85.93	89.53	42.00	99.05
GLIF [57]	111.58	103.07	95.42	95.04	90.15	63.60
LTC [59]	<b>104.10</b>	<u>99.09</u>	86.33	90.94	<b>100.00</b>	<b>100.00</b>
SPSN [73]	120.43	-	83.88	-	45.70	-
TCLIF [54]	286.71	255.67	86.81	92.08	19.10	19.90
LM-H [55]	122.69	102.05	77.70	83.14	99.25	96.10
CLIF [74]	128.28	108.21	43.90	70.44	19.10	64.30
DH-LIF [60]	115.61	<i>100.55</i>	79.12	91.07	98.85	99.35
CELIF [75]	<i>112.35</i>	106.52	<b>97.76</b>	<b>97.66</b>	48.40	<b>100.00</b>
PMSN [58]	113.24	-	<u>96.28</u>	-	<b>100.00</b>	-

incorporates trainable factors to rescale presynaptic inputs at each time step. This technique smooths temporal distributions of gradient norms and stabilizes training. In addition to TDBN and TEBN, we also assess Layer Normalization (LayerNorm) [71], which is commonly used in non-spiking sequence models and can be seamlessly applied to SNNs.

Three major observations can be drawn from the results presented in Table IV. First, TEBN, TDBN, and LayerNorm all significantly enhance performance across the three benchmarks. Second, TEBN, in particular, shows substantial accuracy improvements on the PS-MNIST and binary adding tasks. The large improvement is attributed to the rescaling of presynaptic inputs with trainable factors at each time step in TEBN. This mechanism helps stabilize gradient flow over time and ensures that relevant information is preserved throughout the learning process. Third, LayerNorm proves to be particularly effective in language modeling tasks like PTB. Its ability to normalize across non-batch dimensions facilitates the learning of temporal sequences with varying effective lengths, leading to better performance in handling complex temporal data.

### D. Spiking Neuron Models

In this section, we conduct a comprehensive benchmarking of existing spiking neuron models for temporal processing. The widely used LIF model [61] has limited memory capacity and often suffers from the temporal gradient vanishing problem [54], which results in relatively poor temporal processing capabilities. Recently, several advanced spiking neuron models have been proposed to provide enhanced temporal modeling capabilities. These models can be categorized into two main types. The first category encompasses single-compartment spiking neuron models, which represent biological neurons as undivided units with enriched neuronal dynamics. For instance, Parametric LIF (PLIF) [56] introduces a learnable leaky constant for each neuron, allowing the model to adapt to diverse decaying rates of the input. Adaptive LIF (ALIF) [67] and Liquid Time-Constant (LTC) [59] models incorporate the adaptive threshold as an additional state variable, allowing the spiking neuron

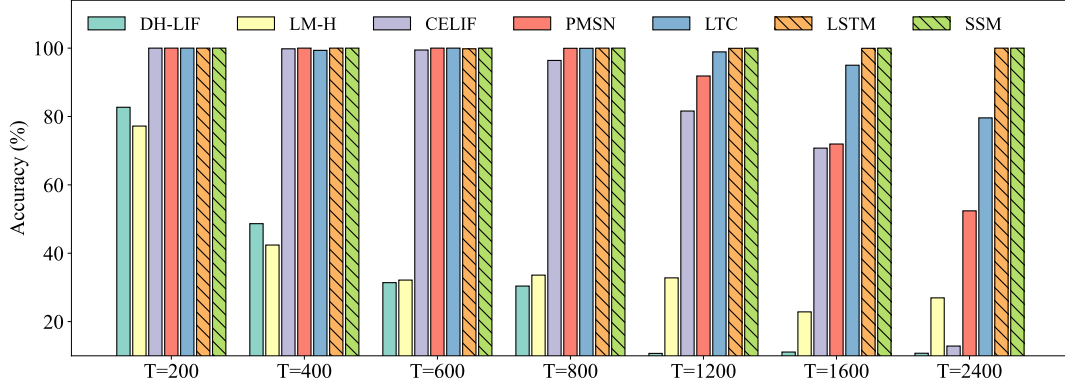


Fig. 8. Comparison of advanced spiking and non-spiking sequence models. Each model is evaluated on the binary adding task with sequence lengths ranging from 200 to 2400.

to retain information related to its firing history within the firing threshold. Furthermore, models such as Generalized Leaky Integrate-and-Fire (GLIF) [57], Complementary Leaky Integrate-and-Fire (CLIF) [74], adaptive Leaky Integrate-and-Fire (adLIF) [53], and Context Embedding Leaky Integrate-and-Fire (CELIF) [75] integrate various biological mechanisms such as ion channels, complementary membrane potentials, adaptation currents, and temporal context to enhance neuronal dynamics of spiking neurons.

In another vein of research, several multi-compartment spiking neuron models have also been proposed. Inspired by P-R neurons in the hippocampus, the Two-Compartment LIF (TC-LIF) model [54] is specifically designed to model interactions between the soma and dendrites. To further reduce the parameter constraints of the TC-LIF model, the Learnable Multi-Hierarchical (LM-H) model [55] has been proposed, demonstrating more stable gradient propagation in deep networks. Going beyond models with meticulously designed two compartments, the Dendritic Heterogeneity LIF (DH-LIF) model [60] endows multiple dendritic compartments with heterogeneous decaying time constants. To address concerns about the slow training speed of SNNs, recent research has also proposed several neuron models to enable parallel computation in the temporal dimension. The family of Parallel Spiking Neural (PSN) models [73] transforms the charging dynamics of the membrane potential into a learnable decay matrix and bypasses the neuron's reset mechanism to enable parallel computation. Furthermore, the Parallel Multi-Compartment Neuron (PMSN) model [58] incorporates multiple interconnected neuronal compartments. These inter-compartment interactions can effectively represent temporal information across diverse timescales.

Since most of these advanced models are developed based on static image datasets, their efficacy for temporal processing tasks remains elusive. To this end, we comprehensively compare their performance against the LIF model on our temporal processing benchmarks. It is important to note that different studies have adopted different numbers of learnable model parameters for their respective neuron models. To ensure a fair comparison, we proportionally adjust the number of neurons in the hidden layers so that all neuron models are evaluated with the same number of parameters as the baseline. As the results

presented in Table V, the various neuron models exhibit varying degrees of improvement over the LIF model, highlighting their effectiveness in enhancing the temporal processing capabilities of SNNs.

The observed improvements can be explained by several specific architectural features of these models. For instance, the improved performance achieved by ALIF can be attributed to the additional slow-decaying state variables, which facilitate information integration over longer timespans. Moreover, the heterogeneous decaying rate of state variables used in these models can further facilitate the establishment of temporal dependencies across different timescales, as demonstrated in models like GLIF and DH-LIF. Furthermore, models such as adLIF, TCLIF, LM-H, and PMSN incorporate recurrent interactions between neuronal variables. These enriched neuronal dynamics significantly enhance the integration of temporal information. Other strategies, such as the extended temporal receptive field employed by SPSN, the temporal context in CELIF, and the liquid decaying rate in LTC, also demonstrate the enhancement of temporal processing capability.

The results also highlight the importance of the time decaying factors in facilitating temporal processing. For instance, the TC-LIF neuron omits these decay factors, leading to continuous accumulation of information in the two membrane potential variables. While this configuration works well for static datasets, it encounters challenges when targets change over time, as observed in tasks such as the PTB and binary adding. This necessitates a neuron model capable of rapidly forgetting past information to respond promptly to new inputs. By mitigating the excessive accumulation of historical inputs, these decay factors ensure that the spiking neuron remains responsive to new stimuli over time.

#### E. Comparison of Spiking and Non-Spiking Sequence Models

Having demonstrated the effectiveness of various advanced SNN models in temporal processing, we further benchmark their performance by comparing them against leading non-spiking sequence models. In the binary adding problem, models such as LM-H, DH-LIF, and PMSN using a feedforward architecture, as well as LTC and CELIF employing a recurrent architecture, can successfully accomplish the task with a sequence length of

TABLE VI

COMPARISON OF DELAY LEARNING AND NEURAL ARCHITECTURES ON TEMPORAL PROCESSING TASKS.  $T_{in}$  REPRESENTS THE INTERNAL TIME WINDOW OF SPIKING NEURONS.

Dataset	PTB ( $T = 70$ )	PS-MNIST ( $T = 784$ )	Binary Adding	
Metric	PPL ↓	Acc. ↑	$T$ ↑	Acc. ↑
#Params.	~5M	~90K	~40K	
LIF	129.96	57.45	100	34.15
LIF w/ DCLS-Delays [76]	89.87	68.98	100	51.85
TCN* [3]	102.20	95.10	1200	69.95
SpikingTCN	114.46	93.76	1200	61.95
LSTM* [2]	88.08	92.41	2400	100
Gated Spiking Neuron [77]	99.98	80.13	1200	29.85
Transformer* [4]	112.43	97.64	2400	100
Spike-Driven Transformer <sup>4</sup> [78]	152.41	96.21	2400	98.15
Spike-Driven Transformer <sup>1</sup> [78]	327.82	95.01	2400	88.05

\* Non-spiking models.

<sup>4,1</sup>  $T_{in} = 4$  and  $T_{in} = 1$ , respectively.

$T = 100$ . These results indicate competitive performance relative to prominent non-spiking RNN models such as Long Short-Term Memory (LSTM) [2] and SSM [5]. To further figure out whether current advanced SNN models already achieve similar to or even surpass the performance of these RNN models, we quantify their long-range temporal processing capacity by varying  $T$  from 100 to 200, 400, 600, 800, 1200, 1600, and 2400.

As shown in Fig. 8, the performance of LM-H and DH-LIF models degrades significantly when  $T$  reaches 400. The performance of CELIF model starts to degrade beyond  $T = 600$  and fails to learn any meaningful temporal information when  $T = 2400$ . Similarly, both LTC and PMSN show varying degrees of degradation as the sequence length increases. In contrast, SSM and LSTM consistently achieve nearly 100% accuracy, even at the challenging scenario when  $T = 2400$ . Collectively, these observations highlight that a significant gap still exists between SNNs and ANNs in modeling long-range dependencies.

#### F. Delay Learning and Neural Architectures

In addition to improving the neuronal dynamics of spiking neurons, many recent studies have also investigated the enhancement of interactions among neurons. This section provides a comprehensive evaluation of these approaches, focusing specifically on the delay learning mechanism and neural architecture designs.

Inspired by neuronal signaling in the brain, delay learning approaches incorporate axonal delay, the time it takes for signals to travel along an axon, into neuron modeling. In our experiments, we utilize a SOTA delay learning method, DCLS-Delays [76], which leverages 1-D temporal convolutions to enable effective delay modeling. The results in Table VI show that the LIF model combined with DCLS-Delays consistently outperforms its counterpart without delay modeling across all three benchmarks. This improvement underscores the significant advantage of delay learning, allowing the model to effectively establish temporal dependencies through the delay line.

In terms of neural architectures, some recent studies focus on adapting advanced non-spiking architectures into spiking counterparts [77], [78]. Here, we compare three well-established neu-

TABLE VII

COMPARISON OF ENERGY EFFICIENCY BETWEEN SPIKING ARCHITECTURES AND THEIR NON-SPIKING COUNTERPARTS ON TEMPORAL PROCESSING TASKS. THE RATIO IS CALCULATED AS THE ENERGY COST OF THE NON-SPIKING ARCHITECTURE DIVIDED BY THAT OF ITS SPIKING COUNTERPART.

Dataset	PTB ( $T = 70$ )		PS-MNIST ( $T = 784$ )		Binary Adding ( $T = 1200$ )	
Energy (nJ)	Cost	Ratio	Cost	Ratio	Cost	Ratio
TCN* [3]	6535.7	-	302.8	-	187.4	-
SpikingTCN	1505.6	4.3	13.4	22.61	6.8	27.56
LSTM* [2]	8244.8	-	393.6	-	158.4	-
Gated Spiking Neuron [77]	1038.7	7.9	21.6	18.2	9.6	16.5
Transformer* [4]	10223.0	-	1420.2	-	1059.8	-
Spike-Driven Transformer <sup>4</sup> [78]	1826.3	5.6	131.3	10.8	109.2	9.7
Spike-Driven Transformer <sup>1</sup> [78]	367.3	27.8	30.0	47.3	23.6	44.9

\* Non-spiking models.

<sup>4,1</sup>  $T_{in} = 4$  and  $T_{in} = 1$ , respectively.

ral architectures for temporal processing: LSTM [2], TCN [3], and Transformer [4], with their spiking variants – Gated Spiking Neuron (GSN) [77], SpikingTCN, and Spike-Driven Transformer [78]. LSTM employs gating mechanisms to dynamically control information storage and forgetting, which effectively alleviates the problems of gradient vanishing and exploding. Similarly, GSN also adopts the gating mechanism to control the storage and forgetting of historical information. Unlike LSTMs, which process temporal sequences through iterative updates, TCNs use dilated convolutions to efficiently capture long-range dependencies. Recently, Transformer architectures have transformed temporal processing through the self-attention mechanism. This architecture demonstrates superior capability in managing long-range temporal dependencies, along with enhanced temporal parallelism and scalability.

For Transformer, its spiking variant replaces the original continuous activation functions with the discrete step function employed in spiking neurons. In addition, an internal time window  $T_{in}$  of the spiking neuron is utilized to calculate the firing rate, thereby expanding the representation space of spiking neurons. However, incorporating this extra time window incurs additional computational overhead. To evaluate its impact, we compare the performance of the Spike-Driven Transformer model with  $T_{in} = 4$  and  $T_{in} = 1$ . To ensure a fair comparison, all architectures are configured with a comparable number of parameters for the same task. For the binary adding task, the sequence length is gradually increased from 100 to 2400 until the architecture no longer performs. Both the maximum sequence length and accuracy are reported in Table VI.

A key observation from Table VI is that, although the performance of spiking architectures is generally inferior to that of their non-spiking counterparts, the performance gap is relatively small, particularly when  $T_{in} = 4$ . While the GSN architecture exhibits a significant accuracy gap compared to LSTM, it outperforms many of the advanced spiking neuron models presented in Table V. Notably, the Spike-Driven Transformer excels in the binary adding task, underscoring its strong capability to model long-range dependencies through self-attention. The relatively lower performance of Spike-Driven Transformer on PTB is likely due to the overfitting with the small dataset size, and it is anticipated that performance would improve on a larger dataset. Furthermore, the Spike-Driven

Transformer experiences a significant drop in accuracy when  $T_{in} = 1$ . This finding underscores the importance of utilizing an additional internal time window to enhance the representational capacity of spiking neurons, which is critical for bridging the performance gap with SOTA non-spiking sequence models.

We further compare the energy efficiency of these spiking architectures with that of their non-spiking counterparts. The details of the energy consumption calculation are provided in Appendix C. As shown in Table VII, spiking architectures are generally an order of magnitude more energy efficient than their non-spiking counterparts. This efficiency stems from the spike-based computation, which relies on efficient Accumulate (AC) operations rather than the more expensive Multiply-Accumulate (MAC) operations typically used in non-spiking architectures. Additionally, Spike-Driven Transformer with  $T_{in} = 4$  consumes approximately four times as much energy as the model with  $T_{in} = 1$ . This indicates that the accuracy improvements afforded by the extra time window come at the expense of significantly increased energy consumption, potentially diminishing the competitiveness of SNNs. This trade-off warrants further investigation. Nonetheless, even with  $T_{in} = 4$ , Spike-Driven Transformer still improves energy efficiency by dozens of times compared to its non-spiking counterpart. These results underscore the significant energy-saving benefit of SNNs, making them particularly advantageous for applications in energy-constrained environments.

## VI. CHALLENGES AND FUTURE DIRECTIONS

In this section, we identify four key challenges in current research on SNNs for temporal processing and propose some future directions to overcome these challenges. The primary challenge arises from the inadequacy of current neuromorphic benchmarks in evaluating SNNs' temporal processing capabilities. This limitation hinders the advancement of SNNs in effectively tackling tasks that involve complex and long-range temporal dependencies. To overcome this challenge, future research should prioritize the development of comprehensive neuromorphic benchmarks specifically designed to evaluate the SNN's abilities to capture long-range temporal dependencies. To ensure practical relevance, these benchmarks should encompass a diverse array of real-world temporal processing tasks that necessitate the maintenance of temporal context over extended durations. Additionally, they should be designed to emphasize the strengths of neuromorphic computing, such as high energy efficiency and low latency. By establishing these benchmarks, we can attain a more accurate assessment of SNN's performance in temporal processing, thereby facilitating their application in real-world temporal signal processing scenarios.

Another major challenge arises from the significant accuracy drops in online learning algorithms when applied to temporal processing tasks. While these algorithms have effectively enhanced the training efficiency and adaptability of SNNs, they often compromise the precision of temporal gradients, leading to suboptimal performance in learning long-range temporal dependencies. Addressing this challenge requires the development of a new generation of online or on-chip learning algorithms that can effectively and efficiently learn

long-range temporal dependencies. This advancement would enable SNNs to continuously adapt to diverse and dynamic real-world scenarios.

The third challenge is that existing SNN models struggle to model long-range dependencies. The primary issues likely arise from optimization challenges, such as the vanishing gradient problem, where gradients associated with earlier time steps become exponentially smaller, resulting in a bias toward short-term dependencies. Additionally, current SNNs face difficulties in effectively storing and retrieving information over time, which hinders SNNs' capabilities in temporal processing. To address this challenge, future research should draw inspiration from the structures and functions of biological neural networks to create spiking neuron models and neural architectures with enhanced memory storage and information retrieval capabilities.

The last challenge concerns the prohibitive training time of SNNs for long sequences. Due to the inherently temporal nature of SNNs, training times in earlier SNN studies scale linearly with sequence length, especially for spiking neuron models that involve non-linear dynamics. This scaling makes it challenging to fully exploit the parallel processing capabilities of hardware accelerators like GPUs. Consequently, training SNNs can be excessively time-consuming, which limits researchers' ability to rapidly validate new ideas and iteratively refine SNN approaches. Despite recent advancements in temporal parallel spiking neuron models that significantly accelerate training, their reliance on linear recurrency as a prerequisite for parallelism fundamentally constrains their ability to capture the rich nonlinear dynamics present in biological neurons. Future work should focus on developing novel spiking neuron models that facilitate parallelized training over time, particularly by supporting the linearization of complex dynamics. Such advancements would drastically reduce training times, thereby enabling rapid development of SNN approaches.

## VII. CONCLUSION

In this work, we introduce the STP, a novel analytical tool designed to evaluate the effectiveness of benchmarks in assessing the temporal processing capabilities of SNN. Our application of the STP to widely used neuromorphic benchmarks indicates that these benchmarks can often be addressed without modeling the temporal dependencies of distant inputs. This finding suggests a significant deficiency in temporal information within current benchmarks, rendering them inadequate for a comprehensive evaluation of various methods' temporal processing capabilities.

To further elucidate the current state of SNN approaches in temporal processing, we developed a suite of benchmarks encompassing three distinct temporal processing tasks and conducted a thorough comparison of existing SNN methodologies. Our benchmarking study demonstrated the enhanced temporal processing capabilities of recently introduced spiking neuron models and architectures. However, it also revealed that, despite the significant energy efficiency of SNNs, there remains a notable performance gap compared to SOTA non-spiking sequence models, particularly in the context of long sequences. Moreover, we discuss the challenges and future directions for SNNs in temporal processing, emphasizing the

urgent need for the development of specialized neuromorphic benchmarks, learning algorithms, and computational models specifically tailored to real-world temporal processing tasks.

## REFERENCES

- [1] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *IEEE transactions on Signal Processing*, vol. 45, no. 11, pp. 2673–2681, 1997.
- [2] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, November 1997.
- [3] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *CoRR*, vol. abs/1803.01271, 2018.
- [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc., 2017.
- [5] A. Gu, K. Goel, and C. Re, "Efficiently modeling long sequences with structured state spaces," in *International Conference on Learning Representations*, 2022.
- [6] L. Deng, G. Li, S. Han, L. Shi, and Y. Xie, "Model compression and hardware acceleration for neural networks: A comprehensive survey," *Proceedings of the IEEE*, vol. 108, no. 4, pp. 485–532, 2020.
- [7] W. Maass, "Networks of spiking neurons: The third generation of neural network models," *Neural Networks*, vol. 10, no. 9, pp. 1659–1671, 1997.
- [8] G. Li, L. Deng, H. Tang, G. Pan, Y. Tian, K. Roy, and W. Maass, "Brain-inspired computing: A systematic survey and future trends," *Proceedings of the IEEE*, vol. 112, no. 6, pp. 544–584, 2024.
- [9] K. Roy, A. Jaiswal, and P. Panda, "Towards spike-based machine intelligence with neuromorphic computing," *Nature*, vol. 575, no. 7784, pp. 607–617, 2019.
- [10] Q. Yu, H. Tang, K. C. Tan, and H. Li, "Rapid feedforward computation by temporal encoding and learning with spiking neurons," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 24, no. 10, pp. 1539–1552, 2013.
- [11] M. Davies, N. Srinivasa, T. Lin, G. N. Chinya, Y. Cao, S. H. Choday, G. D. Dimou, P. Joshi, N. Imam, S. Jain, Y. Liao, C. Lin, A. Lines, R. Liu, D. Mathakutty, S. McCoy, A. Paul, J. Tse, G. Venkataramanan, Y. Weng, A. Wild, Y. Yang, and H. Wang, "Loihi: A neuromorphic manycore processor with on-chip learning," *IEEE Micro*, vol. 38, no. 1, pp. 82–99, 2018.
- [12] J. Pei, L. Deng, S. Song, M. Zhao, Y. Zhang, S. Wu, G. Wang, Z. Zou, Z. Wu, W. He, F. Chen, N. Deng, S. Wu, Y. Wang, Y. Wu, Z. Yang, C. Ma, G. Li, W. Han, H. Li, H. Wu, R. Zhao, Y. Xie, and L. Shi, "Towards artificial general intelligence with hybrid Tianjic chip architecture," *Nature*, vol. 572, no. 7767, pp. 106–111, Aug. 2019.
- [13] D. Ma, J. Shen, Z. Gu, M. Zhang, X. Zhu, X. Xu, Q. Xu, Y. Shen, and G. Pan, "Darwin: A neuromorphic hardware co-processor based on spiking neural networks," *Journal of Systems Architecture*, vol. 77, pp. 43–51, 2017.
- [14] D. Ma, X. Jin, S. Sun, Y. Li, X. Wu, Y. Hu, F. Yang, H. Tang, X. Zhu, P. Lin, and G. Pan, "Darwin3: A large-scale neuromorphic chip with a novel ISA and on-chip learning," *National Science Review*, vol. 11, no. 5, p. nwae102, 03 2024.
- [15] M. Yao, O. Richter, G. Zhao, N. Qiao, Y. Xing, D. Wang, T. Hu, W. Fang, T. Demirci, M. De Marchi, L. Deng, T. Yan, C. Nielsen, S. Sheik, C. Wu, Y. Tian, B. Xu, and G. Li, "Spike-based dynamic computing with asynchronous sensing-computing neuromorphic chip," *Nature Communications*, vol. 15, no. 1, p. 4464, May 2024.
- [16] X. Ju, B. Fang, R. Yan, X. Xu, and H. Tang, "An FPGA implementation of deep spiking neural networks for low-power and fast classification," *Neural Computation*, vol. 32, no. 1, pp. 182–204, 01 2020.
- [17] A. V. Herz, T. Gollisch, C. K. Machens, and D. Jaeger, "Modeling single-neuron dynamics and computations: A balance of detail and abstraction," *Science*, vol. 314, no. 5796, pp. 80–85, 2006.
- [18] J. Wu, E. Yilmaz, M. Zhang, H. Li, and K. C. Tan, "Deep spiking neural networks for large vocabulary automatic speech recognition," *Frontiers in Neuroscience*, vol. 14, p. 199, 2020.
- [19] J. Wu, Y. Chua, M. Zhang, H. Li, and K. C. Tan, "A spiking neural network framework for robust sound classification," *Frontiers in Neuroscience*, vol. 12, p. 836, 2018.
- [20] Q. Yu, R. Yan, H. Tang, K. C. Tan, and H. Li, "A spiking neural network system for robust sequence recognition," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 27, no. 3, pp. 621–635, 2016.
- [21] J. Wu, C. Xu, X. Han, D. Zhou, M. Zhang, H. Li, and K. C. Tan, "Progressive tandem learning for pattern recognition with deep spiking neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 11, pp. 7824–7840, 2022.
- [22] Q. Yang, J. Wu, M. Zhang, Y. Chua, X. Wang, and H. Li, "Training spiking neural networks with local tandem learning," in *Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., 2022, pp. 12 662–12 676.
- [23] Y. Hu, Q. Zheng, X. Jiang, and G. Pan, "Fast-SNN: Fast spiking neural network by converting quantized ANN," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 12, pp. 14 546–14 562, 2023.
- [24] Q. Xu, Y. Li, J. Shen, J. K. Liu, H. Tang, and G. Pan, "Constructing deep spiking neural networks from artificial neural networks with knowledge distillation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2023, pp. 7886–7895.
- [25] S. Deng, Y. Li, S. Zhang, and S. Gu, "Temporal efficient training of spiking neural network via gradient re-weighting," in *International Conference on Learning Representations*, 2022, pp. 1–14.
- [26] Y. Guo, Y. Chen, L. Zhang, X. Liu, Y. Wang, X. Huang, and Z. Ma, "IM-Loss: Information maximization loss for spiking neural networks," in *Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., 2022, pp. 156–166.
- [27] S. Lian, J. Shen, Q. Liu, Z. Wang, R. Yan, and H. Tang, "Learnable surrogate gradient for direct training spiking neural networks," in *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23. International Joint Conferences on Artificial Intelligence Organization*, 8 2023, pp. 3002–3010, main Track.
- [28] Y. Zhu, J. Ding, T. Huang, X. Xie, and Z. Yu, "Online stabilization of spiking neural networks," in *The Twelfth International Conference on Learning Representations*, 2024.
- [29] H. Jiang, G. D. Masi, H. Xiong, and B. Gu, "NDOT: Neuronal dynamics-based online training for spiking neural networks," in *Forty-first International Conference on Machine Learning*, 2024.
- [30] H. Shen, Q. Zheng, H. Wang, and G. Pan, "Rethinking the membrane dynamics and optimization objectives of spiking neural networks," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [31] X. Chen, Q. Yang, J. Wu, H. Li, and K. C. Tan, "A hybrid neural coding approach for pattern recognition with spiking neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 5, pp. 3064–3078, 2024.
- [32] M. Yao, G. Zhao, H. Zhang, Y. Hu, L. Deng, Y. Tian, B. Xu, and G. Li, "Attention spiking neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 8, pp. 9393–9410, 2023.
- [33] Q. Xu, Y. Gao, J. Shen, Y. Li, X. Ran, H. Tang, and G. Pan, "Enhancing adaptive history reserving by spiking convolutional block attention module in recurrent neural networks," in *Advances in Neural Information Processing Systems*, vol. 36. Curran Associates, Inc., 2023, pp. 58 890–58 901.
- [34] Q. Yu, S. Song, C. Ma, L. Pan, and K. C. Tan, "Synaptic learning with augmented spikes," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 3, pp. 1134–1146, 2022.
- [35] Y. Guo, Y. Chen, X. Liu, W. Peng, Y. Zhang, X. Huang, and Z. Ma, "Ternary spike: Learning ternary spikes for spiking neural networks," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 11, pp. 12 244–12 252, Mar. 2024.
- [36] X. Xing, Z. Zhang, Z. Ni, S. Xiao, Y. Ju, S. Fan, Y. Wang, J. Zhang, and G. Li, "SpikeLM: Towards general spike-driven language modeling via elastic bi-spiking mechanisms," in *Forty-first International Conference on Machine Learning*, 2024.
- [37] Y. Hu, Q. Zheng, G. Li, H. Tang, and G. Pan, "Toward large-scale spiking neural networks: A comprehensive survey and future directions," *CoRR*, vol. abs/2409.02111, 2024.
- [38] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [39] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," University of Toronto, Toronto, Ontario, Tech. Rep. 0, 2009.
- [40] J. K. Eshraghian, M. Ward, E. O. Neftci, X. Wang, G. Lenz, G. Dwivedi, M. Bennaamoun, D. S. Jeong, and W. D. Lu, "Training spiking neural networks using lessons from deep learning," *Proceedings of the IEEE*, 2023.
- [41] G. Orchard, A. Jayawant, G. K. Cohen, and N. Thakor, "Converting static image datasets to spiking neuromorphic datasets using saccades," *Frontiers in Neuroscience*, vol. 9, p. 437, 2015.

- [42] H. Li, H. Liu, X. Ji, G. Li, and L. Shi, "CIFAR10-DVS: An event-stream dataset for object classification," *Frontiers in Neuroscience*, vol. 11, p. 244131, 2017.
- [43] A. Amir, B. Taba, D. J. Berg, T. Melano, J. L. McKinstry, C. di Nolfo, T. K. Nayak, A. Andreopoulos, G. Garreau, M. Mendoza, J. Kusnitz, M. DeBole, S. K. Esser, T. Delbrück, M. Flickner, and D. S. Modha, "A low power, fully event-based gesture recognition system," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. IEEE Computer Society, 2017, pp. 7388–7397.
- [44] S. Zhou, B. Yang, M. Yuan, R. Jiang, R. Yan, G. Pan, and H. Tang, "Enhancing SNN-based spatio-temporal learning: A benchmark dataset and cross-modality attention model," *Neural Networks*, vol. 180, p. 106677, 2024.
- [45] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *CoRR*, vol. abs/1804.03209, 2018.
- [46] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, and D. S. Pallett, "DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1," p. 27403, 1993.
- [47] Z. Pan, Y. Chua, J. Wu, M. Zhang, H. Li, and E. Ambikairajah, "An efficient and perceptually motivated auditory neural encoding and decoding algorithm for spiking neural networks," *Frontiers in Neuroscience*, vol. 13, p. 1420, 2020.
- [48] B. Cramer, Y. Stradmann, J. Schemmel, and F. Zenke, "The heidelberg spiking data sets for the systematic evaluation of spiking neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 7, pp. 2744–2757, 2020.
- [49] Y. Wu, L. Deng, G. Li, J. Zhu, and L. Shi, "Spatio-temporal backpropagation for training high-performance spiking neural networks," *Frontiers in Neuroscience*, vol. 12, p. 331, 2018.
- [50] M. Xiao, Q. Meng, Z. Zhang, D. He, and Z. Lin, "Online training through time for spiking neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 20 717–20 730, 2022.
- [51] Q. Meng, M. Xiao, S. Yan, Y. Wang, Z. Lin, and Z.-Q. Luo, "Towards memory- and time-efficient backpropagation for training spiking neural networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 6166–6176.
- [52] E. O. Neftci, H. Mostafa, and F. Zenke, "Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 51–63, 2019.
- [53] A. Bittar and P. N. Garner, "A surrogate gradient spiking baseline for speech command recognition," *Frontiers in Neuroscience*, vol. 16, p. 865897, 2022.
- [54] S. Zhang, Q. Yang, C. Ma, J. Wu, H. Li, and K. C. Tan, "TC-LIF: A two-compartment spiking neuron model for long-term sequential modelling," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 15, 2024, pp. 16 838–16 847.
- [55] Z. Hao, X. Shi, Z. Huang, T. Bu, Z. Yu, and T. Huang, "A progressive training framework for spiking neural networks with learnable multi-hierarchical model," in *The Twelfth International Conference on Learning Representations*, 2023.
- [56] W. Fang, Z. Yu, Y. Chen, T. Masquelier, T. Huang, and Y. Tian, "Incorporating learnable membrane time constant to enhance learning of spiking neural networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 2661–2671.
- [57] X. Yao, F. Li, Z. Mo, and J. Cheng, "GLIF: A unified gated leaky integrate-and-fire neuron for spiking neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 32 160–32 171, 2022.
- [58] X. Chen, J. Wu, C. Ma, Y. Yan, Y. Wu, and K. C. Tan, "PMSN: A parallel multi-compartment spiking neuron for multi-scale temporal processing," *CoRR*, vol. abs/2408.14917, 2024.
- [59] B. Yin, F. Corradi, and S. M. Bohté, "Accurate online training of dynamical spiking neural networks through forward propagation through time," *Nature Machine Intelligence*, vol. 5, no. 5, pp. 518–527, 2023.
- [60] H. Zheng, Z. Zheng, R. Hu, B. Xiao, Y. Wu, F. Yu, X. Liu, G. Li, and L. Deng, "Temporal dendritic heterogeneity incorporated with spiking neural networks for learning multi-timescale dynamics," *Nature Communications*, vol. 15, no. 1, p. 277, 2024.
- [61] A. N. Burkitt, "A review of the integrate-and-fire neuron model: I. homogeneous synaptic input," *Biological Cybernetics*, vol. 95, pp. 1–19, 2006.
- [62] C. Duan, J. Ding, S. Chen, Z. Yu, and T. Huang, "Temporal effective batch normalization in spiking neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 34 377–34 390, 2022.
- [63] Z. Wang, R. Jiang, S. Lian, R. Yan, and H. Tang, "Adaptive smoothing gradient learning for spiking neural networks," in *International Conference on Machine Learning*. PMLR, 2023, pp. 35 798–35 816.
- [64] M. Marcus, B. Santorini, and M. A. Marcinkiewicz, "Building a large annotated corpus of english: The Penn Treebank," *Computational Linguistics*, vol. 19, no. 2, pp. 313–330, 1993.
- [65] G. Bellec, D. Salaj, A. Subramoney, R. Legenstein, and W. Maass, "Long short-term memory and learning-to-learn in networks of spiking neurons," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [66] G. Bellec, F. Scherr, A. Subramoney, E. Hajek, D. Salaj, R. Legenstein, and W. Maass, "A solution to the learning dilemma for recurrent networks of spiking neurons," *Nature Communications*, vol. 11, no. 1, p. 3625, 2020.
- [67] B. Yin, F. Corradi, and S. M. Bohté, "Accurate and efficient time-domain classification with adaptive spiking recurrent neural networks," *Nature Machine Intelligence*, vol. 3, no. 10, pp. 905–913, 2021.
- [68] E. O. Neftci, H. Mostafa, and F. Zenke, "Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 51–63, 2019.
- [69] Y. Bengio, N. Léonard, and A. C. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *CoRR*, vol. abs/1308.3432, 2013.
- [70] H. Zheng, Y. Wu, L. Deng, Y. Hu, and G. Li, "Going deeper with directly-trained larger spiking neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, 2021, pp. 11 062–11 070.
- [71] L. J. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *CoRR*, vol. abs/1607.06450, 2016.
- [72] S. Ioffe and K. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proceedings of the 32nd International Conference on Machine Learning, ICML*, ser. JMLR Workshop and Conference Proceedings, vol. 37. JMLR.org, 2015, pp. 448–456.
- [73] W. Fang, Z. Yu, Z. Zhou, D. Chen, Y. Chen, Z. Ma, T. Masquelier, and Y. Tian, "Parallel spiking neurons with high efficiency and ability to learn long-term dependencies," in *Advances in Neural Information Processing Systems*, vol. 36, 2023, pp. 53 674–53 687.
- [74] Y. Huang, X. Lin, H. Ren, H. Fu, Y. Zhou, Z. Liu, B. Pan, and B. Cheng, "CLIF: Complementary leaky integrate-and-fire neuron for spiking neural networks," in *Proceedings of the 41st International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 235. PMLR, 21–27 Jul 2024, pp. 19 949–19 972.
- [75] X. Chen, J. Wu, H. Tang, Q. Ren, and K. C. Tan, "Unleashing the potential of spiking neural networks for sequential modeling with contextual embedding," *CoRR*, vol. abs/2308.15150, 2023.
- [76] I. Hammouamri, I. K. Hassani, and T. Masquelier, "Learning delays in spiking neural networks using dilated convolutions with learnable spacings," in *The Twelfth International Conference on Learning Representations, ICLR*. OpenReview.net, 2024.
- [77] X. Hao, C. Ma, Q. Yang, J. Wu, and K. C. Tan, "Towards ultra-low-power neuromorphic speech enhancement with spiking-fullsubnet," *CoRR*, vol. abs/2410.04785, 2024.
- [78] M. Yao, J. Hu, Z. Zhou, L. Yuan, Y. Tian, B. Xu, and G. Li, "Spike-driven transformer," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [79] W. Fang, Y. Chen, J. Ding, Z. Yu, T. Masquelier, D. Chen, L. Huang, H. Zhou, G. Li, and Y. Tian, "SpikingJelly: An open-source machine learning infrastructure platform for spike-based intelligence," *Science Advances*, vol. 9, no. 40, p. eadi1480, 2023.
- [80] B. Yin, F. Corradi, and S. M. Bohté, "Accurate and efficient time-domain classification with adaptive spiking recurrent neural networks," *Nature Machine Intelligence*, vol. 3, no. 10, pp. 905–913, Oct 2021.
- [81] I. Loshchilov and F. Hutter, "Fixing weight decay regularization in adam," *CoRR*, vol. abs/1711.05101, 2017.
- [82] S. Merity, N. S. Keskar, and R. Socher, "Regularizing and optimizing LSTM language models," in *6th International Conference on Learning Representations, ICLR*. OpenReview.net, 2018.
- [83] M. Horowitz, "1.1 computing's energy problem (and what we can do about it)," in *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*. IEEE, 2014, pp. 10–14.
- [84] Y. Li, J. Li, K. Sun, L. Leng, and R. Cheng, "Towards scalable GPU-accelerated SNN training via temporal fusion," in *Artificial Neural Networks and Machine Learning – ICANN*. Cham: Springer Nature Switzerland, 2024, pp. 58–73.

# Appendix

## “Spiking Neural Networks for Temporal Processing: Status Quo and Future Prospects”

Chenxiang Ma\*, Xinyi Chen\*, Yanchen Li\*, Qu Yang, Yujie Wu, Guoqi Li, *Member, IEEE*, Gang Pan, *Senior Member, IEEE*, Huajin Tang, *Senior Member, IEEE*, Kay Chen Tan, *Fellow, IEEE*, Jibin Wu, *Member, IEEE*

This appendix provides the implementation details, the organization of our source code, and supplementary visualizations for the analysis of the DvsGesture dataset.

### APPENDIX A

#### IMPLEMENTATION DETAILS OF SECTION III “EVALUATION OF NEUROMORPHIC BENCHMARKS USING STP”

##### A. Datasets

We analyze the effectiveness of ten widely used neuromorphic benchmarks in the evaluation of temporal processing capabilities, including three static image recognition datasets (i.e., MNIST [38], CIFAR10 [39], and CIFAR100 [39]), three event-based vision datasets (i.e., N-MNIST [41], DVS-CIFAR10 [42], and DvsGesture [43]), and four audio classification datasets (i.e., GSC [45], SHD [48], SSC [48], and TIMIT [46]). Details of these datasets and the data augmentation techniques adopted in our experiments are provided below.

- **MNIST** dataset consists of 60,000 training and 10,000 testing handwritten-digit images in 10 classes. The size of each image is  $28 \times 28$ .
- **CIFAR10** and **CIFAR100** datasets contain a total of 50,000 training images and 10,000 test images, categorized into 10 and 100 classes, respectively. For the training set, we apply standard data augmentation, which includes padding each sample with 4 pixels on all sides, followed by a  $32 \times 32$  crop and a random horizontal flip. In line with previous studies [63], we also utilize the AutoAugment and Cutout techniques for additional data augmentation.
- **N-MNIST** dataset is created by capturing static MNIST images with a DVS camera. Each spike pattern has a spatial dimension of  $34 \times 34 \times 2$  and lasts for 300 time steps. No data augmentation methods are used.
- **CIFAR10-DVS** dataset is obtained from the CIFAR-10 dataset by scanning each image with repeated closed-loop movements in front of a DVS camera. It includes 9,000 training samples and 1,000 testing samples, with a spatial resolution of  $128 \times 128$ . Like CIFAR-10, CIFAR10-DVS comprises 10 classes. We utilize the standard preprocessing pipeline from SpikingJelly [79] to convert events into frames for further analysis, without applying any data augmentation techniques.
- **DvsGesture** dataset comprises 11 different hand gestures performed by 29 different subjects, captured using a DVS camera under three varying illumination conditions. We apply the standard preprocessing pipeline from SpikingJelly to transform events into frames with 20 time steps, without implementing any data augmentation techniques.
- **GSCv2** dataset consists of 105,829 one-second audio clips featuring 35 different spoken commands, which are recorded by various speakers in diverse environments. Following existing studies [54], [67], the original classes are restructured into 12 categories, including ten words: “yes”, “no”, “up”, “down”, “left”, “right”, “on”, “off”, “stop”, and “go”, along with a special “unknown” class that encompasses the remaining commands, plus an extra “silence” class extracted from background noise. The preprocessing methods align with those used in prior studies [54], [67].
- **SHD** dataset is a keyword spotting task for spike-based audio classification, featuring spoken digits from 0 to 9 in both English and German, which are categorized into 20 classes. The dataset comprises recordings from twelve distinct speakers, with two exclusively in the test set. Based on the method from previous work [54], each original waveform has been transformed into spike trains across 700 input channels. The training set includes 8,332 samples, while the test set contains 2,088 samples, with no separate validation set.
- **SSC** dataset is created from the GSCv2 dataset, utilizing a biologically inspired encoding method to represent data in a spike-based format. It encompasses 35 classes contributed by a diverse array of speakers. In accordance with the method outlined by Zhang *et al.* [54], the original waveforms have been converted into spike trains across 700 input channels. The dataset includes 75,466 samples in the training set and 20,382 samples in the test set.
- **TIMIT** dataset is a corpus of read speech that includes time-aligned orthographic, phonetic, and word transcriptions, as well as a single-channel, 16-bit, 16 kHz speech waveform file for each utterance. It features broadband recordings from 630 speakers — approximately 70% male and 30% female — representing 8 major dialects of American English, with each speaker reading 10 phonetically rich sentences. Following previous work [80], the raw audio data is preprocessed into Mel Frequency Cepstral Coefficients and converted into frames with 39 input channels. Each frame is then classified into one of 61 phoneme classes for prediction.

##### B. Training Configurations

Our training configurations are detailed in Table VIII. For clarity, the following notations are employed in denoting network architectures: C represents a convolutional layer, AP stands for average pooling, and FC denotes a fully connected layer. Furthermore, C3 indicates a convolutional of kernel size  $3 \times 3$ . The numerical value preceding C and FC indicates the number of output channels.

### APPENDIX B

#### IMPLEMENTATION DETAILS OF SECTION IV “TEMPORAL PROCESSING BENCHMARK SUITE”

The training configurations are outlined in Table VIII. Specifically, for the PS-MNIST and Binary Adding tasks, we utilize the AdamW optimizer [81] in conjunction with a step scheduler that reduces the learning rate by a factor of 0.8 every 10 epochs. In contrast, the PTB task employs the stochastic gradient descent (SGD) optimizer without a scheduler. This configuration has been validated as effective for training the SNN methods used in this study. Fully connected network architectures are implemented for all SNNs in the experiments. Consistent with standard practices, the hidden dimensions are set to 64FC–256FC–256FC for the PS-MNIST task [54], [67] and 400FC–1100FC for the PTB task [82]. For our proposed Binary Adding task, the hidden dimensions are configured as 128FC–128FC.

TABLE VIII  
TRAINING CONFIGURATIONS AND HYPER-PARAMETERS FOR THE EVALUATION OF NEUROMORPHIC BENCHMARKS

Dataset	Epochs	Optimizer	Learning Rate	Learning Rate Schedule	Batch Size	Time Step ( $T$ )	Neuronal Decay	Threshold	Network Architecture
MNIST	100	AdamW	0.0001	Cosine Annealing	256	10	0.5	0.3	32C3-AP2-32C3-AP2-128FC-10FC
CIFAR10	200	SGD	0.1	Cosine Annealing	128	4	0.3	1.0	ResNet18
CIFAR100	200	SGD	0.1	Cosine Annealing	128	4	0.3	1.0	ResNet18
N-MNIST	100	AdamW	0.0001	Cosine Annealing	16	300	0.65	0.3	64C7-AP2-128C7-128C7-AP2-10FC
DVS-CIFAR10	200	AdamW	0.001	Cosine Annealing	64	10	0.3	1.0	VGG11
DvsGesture	200	AdamW	0.001	Cosine Annealing	32	20	0.3	1.0	VGG11
GSC	200	AdamW	0.01	Cosine Annealing	128	101	0.8	0.5	(512FC)*6-10FC
SHD	200	AdamW	0.0005	-	256	250	0.9	0.7	(128FC)*5-20FC
SSC	200	AdamW	0.0005	-	256	250	0.9	0.9	(128FC)*5-35FC
TIMIT	200	AdamW	0.0005	-	64	100	0.7	0.7	(1024FC)*5-61FC
PTB	100	SGD	3	-	20	70	0.5	0.5	400FC-1100FC-400FC
PS-MNIST	100	AdamW	0.0005	StepLR	256	784	1.0	0.5	64FC-256FC-256FC-10FC
Binary Adding	50	AdamW	0.0005	StepLR	250	100	0.98	0.5	128FC-128FC-10FC

## APPENDIX C

### DETAILS OF ENERGY COST COMPUTATION FOR NON-SPIKING AND SPIKING NEURAL ARCHITECTURES

This section outlines the computation of empirical energy cost metrics presented in Table VII. We first establish theoretical energy cost formulas for each hidden layer of compared models, as presented in Table IX. Following standard computational cost evaluation approaches for ANNs and SNNs, we count their number of 32-bit floating-point Multiply-Accumulate (MAC) and Accumulate (AC) operations per inference based on their spatiotemporal dynamics. These numbers of operations are then used to estimate the energy consumption of models on neuromorphic hardware, where  $E_{MAC}$  and  $E_{AC}$  denote the energy required per MAC and AC operation, respectively. In these formulas,  $m$  and  $n$  denote the input sizes, and hidden sizes for each layer, respectively. The variables  $f_{in}$  and  $f_{out}$  represent the spike frequencies of input and output spike sequences. For the SpikingTCN,  $k$  denotes the convolution kernel size,  $f_{conv2}$  specifies the spike frequency of input spikes in the second convolution layer. In the Spike-Driven Transformer,  $h$  is the hidden dimension of feedforward modules,  $f_Q$ ,  $f_K$ ,  $f_V$ ,  $f_{attn}$ ,  $f_{fc1}$ , and  $f_{fc2}$  denote the spike frequency of neurons in the self-attention and feedforward modules. These spike frequency data is recorded during inference across three tasks and used to compute the number of AC operations.

Subsequently, based on data derived from a 45 nm CMOS process [83], the energy per operation is measured to  $E_{AC} = 0.9$  pJ for AC operations and  $E_{MAC} = 4.6$  pJ for MAC operations. All the above data are ultimately substituted into the theoretical energy cost formulas, computing the empirical energy cost for each layer. The total energy cost is then obtained by summing the energy consumption across all hidden layers, as demonstrated in the main text.

TABLE IX  
THEORETICAL ENERGY COSTS OF DIFFERENT NEURAL ARCHITECTURES

Architecture	Theoretical Energy Cost
TCN	$(kmn + kn^2) \cdot E_{MAC}$
SpikingTCN	$(kmn \cdot f_{in} + kn^2 \cdot f_{conv2}) \cdot E_{AC}$
LSTM	$(4mn + 4n^2 + 19n) \cdot E_{MAC}$
GSU	$(2mn \cdot f_{in} + 2n^2 \cdot f_{out}) \cdot E_{AC} + 5n \cdot E_{MAC}$
Transformer	$(4n^2 + 2nT + 2nh) \cdot E_{MAC}$
Spike-Driven Transformer ( $T_{in} = 4$ )	$((12f_{in} + 4f_{attn}) \cdot n^2 + (4f_Q \cdot f_K + 4f_V) \cdot nT + (4f_{fc1} + 4f_{fc2}) \cdot nh) \cdot E_{AC} + (24n + 4h) \cdot E_{MAC}$
Spike-Driven Transformer ( $T_{in} = 1$ )	$((3f_{in} + f_{attn}) \cdot n^2 + (f_Q \cdot f_K + f_V) \cdot nT + (f_{fc1} + f_{fc2}) \cdot nh) \cdot E_{AC}$

## APPENDIX D

### ORGANIZATION OF THE BENCHMARKING LIBRARY

In this section, we provide a detailed description of the structure of our open-sourced benchmarking library. Initially, we present an overview of the library, highlighting its main components such as the framework and experimental configurations. We then examine the structure of the model design within the framework, elucidating the role and functionality of each component. This systematic exposition is designed to facilitate the rigorous evaluation and further development of SNN models and datasets in future studies.

#### A. Overview of the Benchmarking Library

This subsection presents the architecture of our benchmarking library, which comprises two key components: **framework** and **experiments**. Together, these components establish a comprehensive structure that supports both model development and task design while ensuring a well-

TABLE X  
OVERVIEW OF THE BENCHMARKING LIBRARY

Module	Component	Instance	Description
kernel	-	temporal_fusion_kernel, etc.	Ready-to-use accelerated CUDA kernel wrappers optimized for spiking neurons.
network	neuron structure trainer	LIF, ALIF, TDBN, TEBN, etc. GSU, SpikingTCN, etc. SurrogateGradient, etc.	Interfaces for modules defining network layers, architectures, and the corresponding training methods required for SNNs.
utils	dataset tools	PTB, PS-MNIST, etc. logging, save_checkpoint, etc.	Encapsulation of utility functions and optimized dataset modules for efficient model training and validation.

defined experimental setup. They collectively provide essential interfaces that enable users to interact with the benchmark and contribute new features.

The **framework** serves as the foundation for model design, offering functionalities for defining SNN models and selecting appropriate training strategies. This component plays a crucial role in deploying models for specific tasks and incorporates Compute Unified Device Architecture (CUDA)-based acceleration kernels along with preprocessing capabilities tailored for task-specific datasets.

The **experiments** component provides a detailed set of benchmark experiments, serving as practical references for users. Each experiment is organized as an independent module, consisting of a code file and a configuration file. The code file specifies the experimental workflow, while the configuration file simplifies the process by enabling users to load and manage experiment-specific configurations, including hyperparameter settings efficiently.

### B. Structure of framework

The **framework** component of the benchmarking library is designed to assist users in defining, training, and deploying SNN models through structured interfaces. It consists of three key modules: **kernel**, **network**, and **utils**. A detailed breakdown of these components and their specific functionalities is presented in Table X. The **kernel** module facilitates CUDA-based acceleration to enhance computational efficiency, incorporating an optional kernel based on the temporal fusion method [84], which is specifically designed to accelerate processing within spiking neuron layers during long temporal sequences in both training and inference.

Within the **network** module, three primary components are provided: **neuron**, **structure**, and **trainer**. The **neuron** component enables users to select particular spiking neurons. The **structure** component offers a collection of predefined network architectures. Once the model structure is established, the **trainer** component provides various training strategies to optimize model performance.

The **utils** module offers essential functionalities for both training and deployment. It includes the **tools** component, which provides optional utilities to facilitate model training, and the **dataset** component, which supports various neuromorphic datasets.

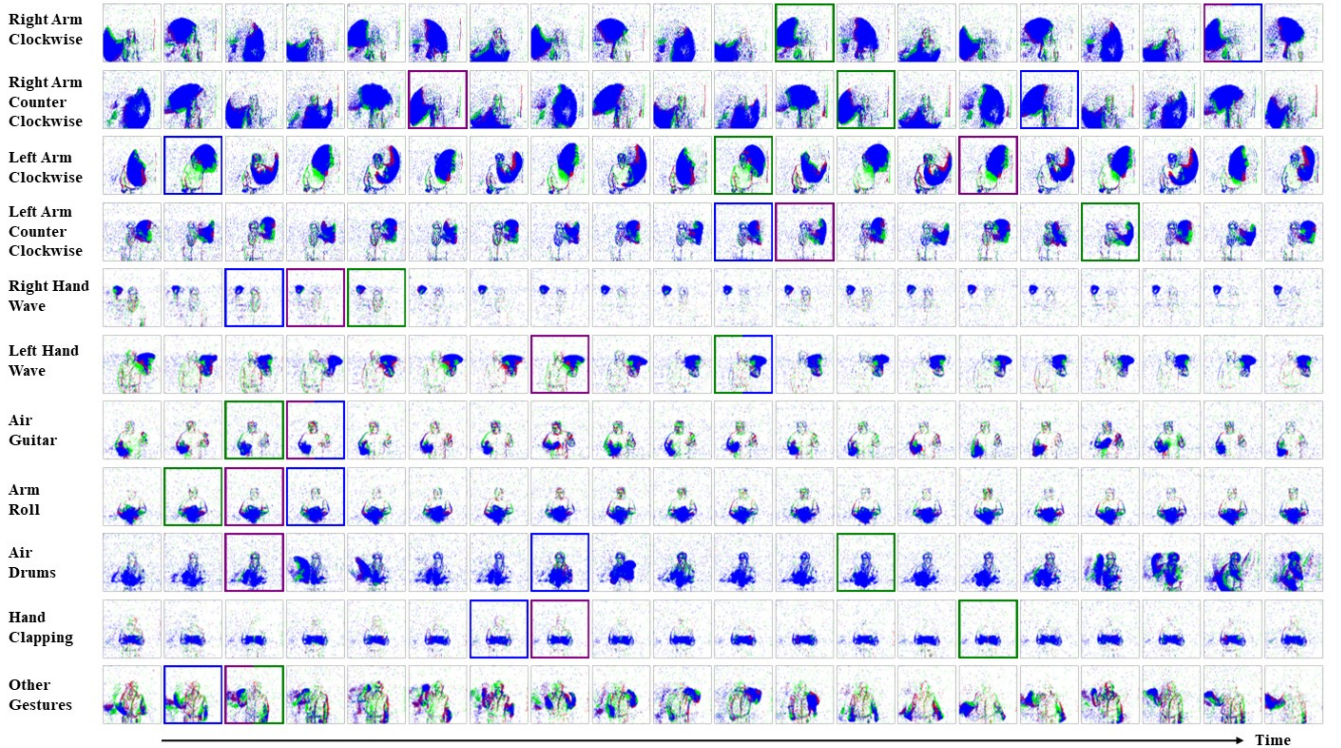


Fig. 9. Qualitative results of samples from the DvsGesture dataset. In these samples, the most confident frame differs across the three algorithms in STP. Frames selected by STBP are highlighted with purple boxes, SDBP with green boxes, and NoTD with blue boxes.

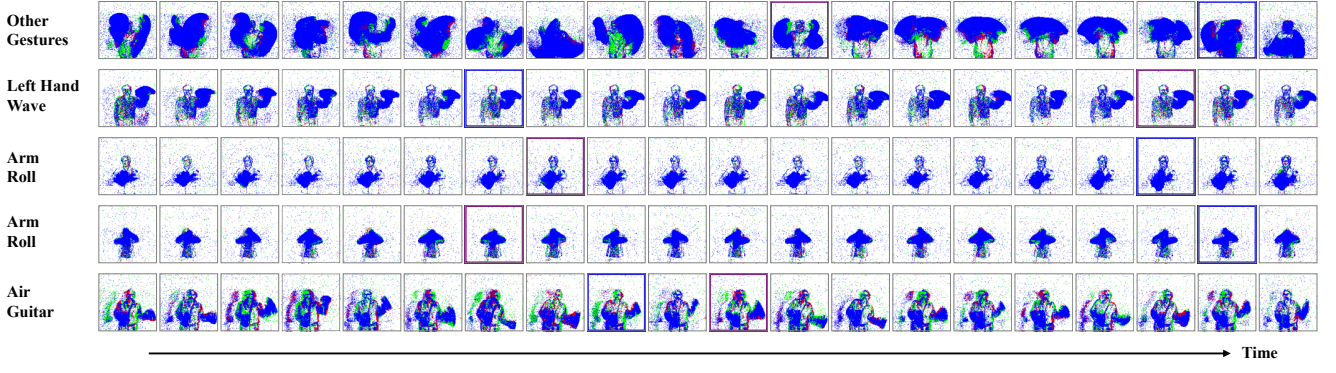


Fig. 10. Visualization of samples from the DvsGesture dataset that STBP classifies correctly, while NoTD does not. The predictions by NoTD are right arm counter clockwise, left arm clockwise, air drums, hand clapping, and other gestures (from top to bottom).

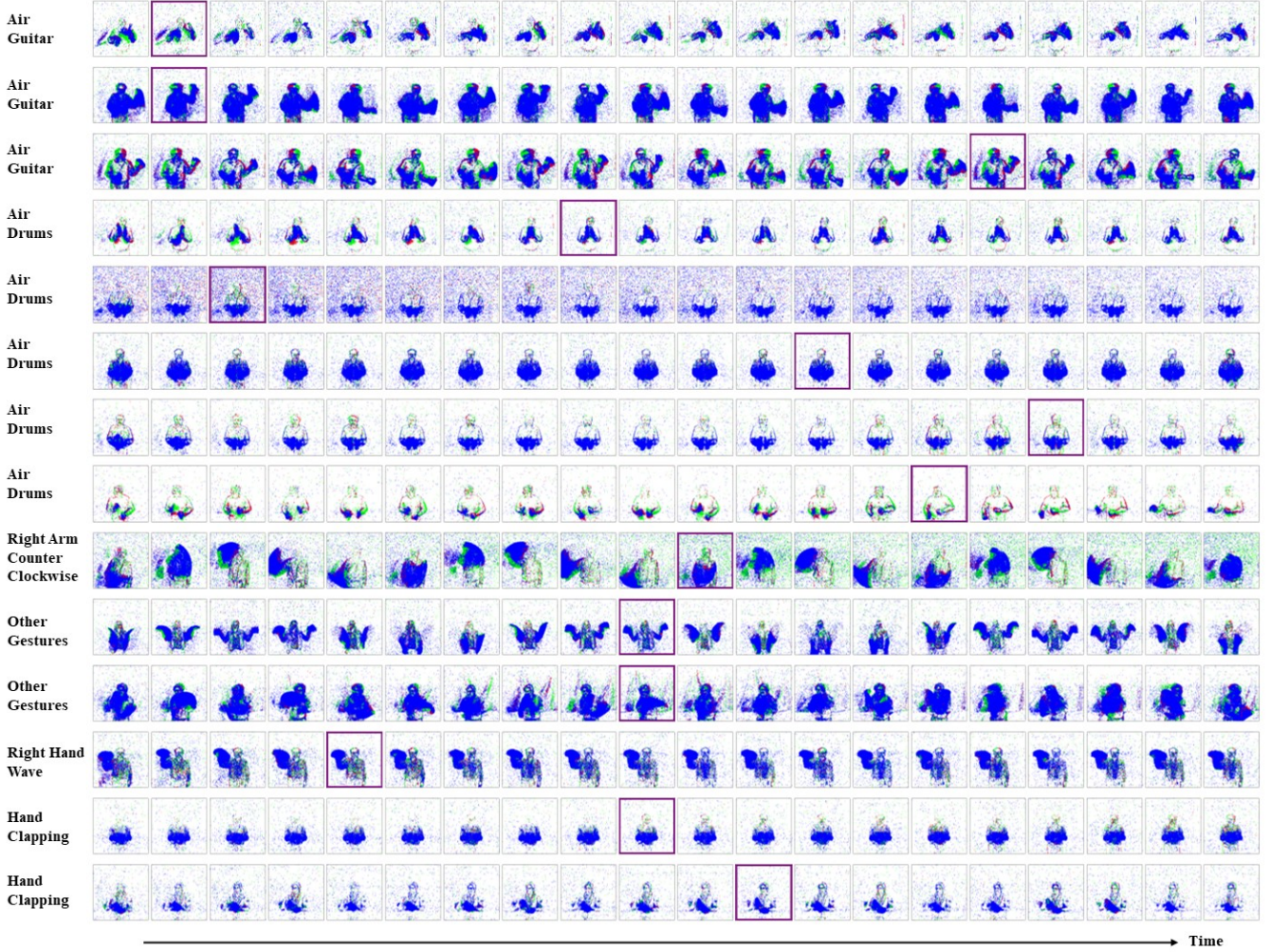


Fig. 11. Visualization of samples from the DvsGesture dataset that are misclassified by STBP. The predictions by STBP are other gestures, air drums, other gestures, hand clapping, arm roll, arm roll, hand clapping, other gestures, other gestures, air guitar, air drums, right arm counter clockwise, arm roll, air drums (from top to bottom).